

5

## METHODS AND SYSTEM FOR THE IDENTIFICATION AND CHARACTERIZATION OF PEPTIDES AND THEIR FUNCTIONAL RELATIONSHIPS BY USE OF MEASURES OF CORRELATION

10

### Technical Field of the Invention

The invention relates to the field of high-throughput analysis of samples with peptide content and especially computer implemented methods and a system implementing these methods for identifying and characterizing peptides and their functional relationships by use of measures of correlation.

15

### Background of the Invention

The success of the Human Genome Project in mapping the human genetic code offers astonishing potential for medical research. A prerequisite for using this information, however, is the identification of gene products, especially proteins and peptides. Peptides are the family of molecules formed from the linking, in a defined order, of various amino acids. The link between one amino acid residue and the next is an amide bond, and is sometimes referred to as a peptide bond. Peptides occur in nature and are responsible for a variety of functions, many of which are not understood. They differ from proteins, which are also long chains of amino acids, by virtue of their size.

25

Parallel to the world-wide efforts in Genomics, a variety of discovery technologies have been developed for analyzing samples with peptide content. Just as Genomics focuses on decoding the human genome, these technologies strive for an comprehensive analysis of the myriad of biologically relevant proteins and peptides with a molecular mass between about 0.5 and 20 kDa, among which insulin is a prominent example.

30

Profiling of peptides and proteins of human body fluids and tissues by mass spectrometry reveals a large number of peptide signals. Such high-throughput analytical processes demand highly sophisticated bioinformatic approaches to understand and analyze biological and pharmaceutical coherences in huge sets of data.

35

Conventional computer implemented methods for assisting the mass spectrometric identification of peptides and small proteins interpret the spectra and generate proposals for

the identity of the candidate peptide signal by determining the differences of masses of the fragments in one spectra and assigning these differences to missing amino acids. A string of missing amino acids is then composed to a proposed amino acid sequence that is thereafter queried in a huge database containing tens of thousands of the known protein sequences, 5 such as the Swiss-Prot database. However, if the analyzed peptide or protein is not abundant and/or in a complex mixture, such an approach turns out to be not very effective and, thus, time consuming concentration or fractionation steps of the sample have to be performed.

10 More sophisticated approaches take the knowledge of a known sequence in a spectrum into consideration. Here, proteolytic digests of the known sequence are proposed "in silico", and a hypothetic resulting spectrum is then correlated with the actually measured one. However, these approaches are successful only if the sources of the spectra contain only a few different analytes, as their fragment signals alter the calculations and lower the correlation 15 coefficients of the hypothesized calculated spectra with the actually measured one. If many possible protein precursors exist for a given peptide, then creating such a hypothetic spectrum for each unknown peptide and each possible precursor, the correlation process of hypothetic and measured spectra often turn out to be quite laborious and at times even unsuccessful.

20 Eng et al (*Journ. Am. Soc. Mass Spectrom.* 5, 976-989, 1994) for instance describe a statistical scorer for tandem mass spectrometry, that relies on cross-correlating experimental spectra with predicted spectra of peptides from a database (Havilio et al, *Anal. Chem.* 75[3], 435-444, 2003). No additional information (e.g. enzyme specificity used to create the 25 peptides) about the peptide except the mass of the peptide is used. In a first step the tandem mass spectrometry data is reduced, whereby all but the most abundant signals are removed. In a second step protein sequences are queried from a database for combinations of amino acids that match the mass of the peptide, wherein the search algorithm only considers mass changes typical for a post-translational modification at every occurrence of 30 the modification site. In a third step, the preliminary matches are scored by summing the number of fragmented ions that match the ions observed in the spectrum. Immonium ions are considered if the sequence contains the amino acids Tyrosine, Tryptophan, Methionine or Phenylalanine. This and the sum of fragments are being taken into account in the scoring function. Finally, a spectrum is reconstructed from the putative amino acid sequences and 35 the highest scoring predictions are assessed by a cross-correlational analysis. The cross-correlation function measures the coherence of the reconstructed and the measured spectrum signals by, in effect, translating one signal across the other. Well-known

applications such as SEQUEST and Sonar make use of this approach. However, a disadvantage of this approach is that the peak intensity strongly depends on the ion type, the ion mass and other experimental parameters and that many factors are not fully understood yet that contribute to peptide fragmentation.

5

Perkins et al (*Electrophoresis* 20[18], 3551-3567, 1999) describe a statistical scorer that evaluates the probability of finding a collection of detected fragments in a protein database (Havilio et al, *Anal. Chem.* 75[3], 435-444, 2003). Applications such as Mascot, MOWSE, Protocall are based on this approach. However, a disadvantage of this approach is that the 10 signal intensities of the measured spectra are not being considered for the data analysis.

Weinberger et al (United States Patent Application 2002/0182649) describe essentially two approaches. In the first approach a protein candidate is identified by providing the mass spectrum to a protein database mining tool which identifies at least one protein candidate for 15 the test protein in the database based on a closeness-of-fit measure between the mass spectrum and the theoretically calculated mass spectra of proteins in the database. In the second approach, the protein candidate is directly sequenced using mass spectrometry. In this method the unknown peptide is directly fragmented during mass spectrometry and the masses of the generated fragments are determined by mass spectrometry and are used to 20 calculate the sequence of the unknown peptide.

The approaches according to Eng et al and Weinberger et al have in common that a closeness-of-fit analysis or a cross-correlation is performed over all signals of two spectra, i.e. the measured spectrum and the predicted spectrum. A fundamental disadvantage of 25 these methods is that they rely on predicted mass spectra.

Thus, all of the above approaches have their disadvantages in that at times they turn out to be not very effective, quite laborious, time consuming and often unsuccessful.

30 There is thus a need for methods for analyzing samples with peptide content and a system implementing these methods overcoming or at least mitigating the disadvantages associated with the prior art.

#### Summary of the Invention

35 The following methods according to the present invention are based upon the concept of Correlation Associated Networks and peptide topologies as will be apparent from the detailed description in the sections further below.

According to the present invention a method based on CANs is provided for providing a representative, non-redundant overview of the peptide content of a sample type by analyzing a plurality of samples using its peptide topology, wherein the method comprises the steps of providing a respective mass spectrum for each sample of said plurality of samples, wherein  
5 signal intensity peaks correspond to potential peptides, computing the measures of correlation between the signal intensities of said potential peptides, grouping potential peptides together, which exhibit a degree of correlation among each other above a certain threshold, thereby providing a plurality of correlation associated networks of potential peptides, and assigning one representative potential peptide out of each correlation  
10 associated network as a representative peptide to said correlation associated network of said sample type.

Furthermore, a method based on CANs is provided for predicting the sequence of peptides using the peptide topology of a plurality of samples containing a peptide having a known  
15 precursor, wherein the method comprises the steps of providing a respective mass spectrum for each sample of said plurality of samples, wherein signal intensity peaks correspond to potential peptides, identifying said peptide having a known precursor using the mass of said peptide, wherein the sequence of the known precursor is known, computing the measures of correlation between the signal intensity of said peptide having a known precursor and the  
20 signal intensities of the other potential peptides, selecting potential peptides, which exhibit a degree of correlation with said peptide having a known precursor above a certain threshold, and predicting the sequence of the potential peptides by matching masses of putative fragments of the sequence of the known precursor with the masses of the potential peptides correlating with said peptide having a known precursor.

25 Still furthermore, a method based on CANs is provided for predicting the sequence of peptides using the peptide topology of a plurality of samples containing a peptide with a known sequence, wherein the method comprises the steps of providing a respective mass spectrum for each sample of said plurality of samples, wherein signal intensity peaks correspond to potential peptides, identifying a peptide with a known sequence using its mass, computing the measures of correlation between the signal intensity of said known peptide and the signal intensities of the potential peptides, selecting potential peptides, which exhibit a degree of correlation with the known peptide above a certain threshold, computing the mass differences between each of the potential peptides and the known  
30 peptide, and predicting the sequence and/or the biologically, chemically or physically modified sequence of the potential peptides by using data about mass differences caused by  
35

biological, chemical or physical processes matching the mass differences determined in the previous step.

Yet still furthermore, a method based on CANs is provided for identifying peptides suitable to  
5 be used as marker panels using the peptide topology of a plurality of samples taken from at least two different experimental groups representing a status A and a status B, wherein the method comprises the steps of providing a respective mass spectrum for each sample of said plurality of samples, wherein signal intensity peaks correspond to potential peptides, computing the measures of correlation between the signal intensities of said potential peptides for each plurality of samples within each experimental group separately, and  
10 selecting pairs of potential peptides, which exhibit a difference in the degree of correlation between the different experimental groups above a certain threshold, thereby providing peptides which are suitable to be used as marker panels for diagnostic purposes to distinguish between status A and status B.

15

Yet still furthermore, a method based on CANs is provided for identifying peptides suitable to be used as marker panels using the peptide topology of a plurality of samples taken from at least two different experimental groups representing a status A and a status B, wherein the method comprises the steps of providing a respective mass spectrum for each sample of  
20 said plurality of samples, wherein signal intensity peaks correspond to potential peptides, selecting potential peptides correlating with a parameter being representative of status A or status B, computing the measures of correlation between the signal intensities of said selected potential peptides for each plurality of samples, and selecting pairs of potential peptides which exhibit no correlation of their respective signal intensities above a certain  
25 threshold, thereby providing potential peptides which are suitable to be used as complementing peptides in a marker panel for diagnostic purposes to distinguish between status A and status B.

Finally, a method based on CANs is provided for identifying peptides suitable as a surrogate  
30 for a known peptide using the peptide topology of a plurality of samples, wherein the method comprises the steps of providing a respective mass spectrum for each sample of said plurality of samples, wherein signal intensity peaks correspond to potential peptides, computing the measures of correlation between the signal intensity of said known peptide and the signal intensities of potential peptides, and selecting potential peptides, which exhibit  
35 a degree of correlation with said known peptide above a certain threshold, thereby providing potential peptides suitable as a surrogate for said known peptide.

Preferred embodiments of the present invention are disclosed in the dependent claims.

#### Brief Description of the Drawings

FIGURE 1 shows schematically the hardware components and software modules according to the present invention, their interfaces as well as the flow of information between the hardware components and the software modules.

FIGURE 2 shows an averaged Peptide Mass Fingerprint of cerebrospinal fluid (CSF) samples from several patients. Each of the 96 chromatographic fractions of each sample is analyzed by MALDI-ToF-mass spectrometry and all 96 mass spectra generated from one sample are visualized as a "2-D gel-like picture", wherein the x- and y-axis correspond to the mass-to-charge ratio ( $m/z$ ) and the chromatographic fraction (F), respectively. The bars represent peptide peaks, wherein the color intensity represents the mass spectrometric signal intensity. Some identified peptides including amino acid numbers are identified in this map.

FIGURE 3 shows a diagram exemplifying the correlational behavior of functionally related peptides. Four traces of spectra from four different samples are focused on the signals of a human osteopontin peptide being comprised of the amino acids 249-314 of human osteopontin ( $m/z$  = 7653.6 Da) and its phosphorylated derivatives, carrying one ( $m/z$  = 7733.5 Da), two ( $m/z$  = 7813.5 Da), three ( $m/z$  = 7893.4 Da) or more phosphorylated residues. The conserved concentration ratios of the peptides between samples lead to high degrees of correlation of the signal intensities of the respective peptide pairs.

FIGURE 4 shows a schematic example of a correlation associated network (CAN) according to the present invention. Any CAN starts from a hub peptide and any member of a 1<sup>st</sup> order neighborhood of such a hub peptide of 1<sup>st</sup> order can be a hub peptide for the next order neighborhood and so forth.

FIGURE 5 shows a flow chart schematizing the procedural steps of an application of the CAN Module according to the present invention.

FIGURE 6 shows a graphical representation of an exemplary peptide topology of a sample, wherein peptides are represented by bullets and their mutual relations by lines connecting these bullets. Such a peptide network can be projected onto a peptide map like Figure 2 for a more intuitive analysis of the results.

FIGURE 7 shows a flow chart schematizing the procedural steps of an interaction of the Sequence Network Module with the CAN Module according to the present invention.

5 FIGURE 8a shows a flow chart schematizing the process of checking whether a predicted sequence matches the experimental properties of an unknown peptide.

FIGURE 8b shows a flow chart exemplifying the generation of sequence predictions, which are checked according to Figure 8a.

10 FIGURE 8c shows a flow chart schematizing the query of all unknown peptides P2 which are related to a known peptide P1. Sequence predictions are generated for any unknown peptide P2 according to Figure 8b.

15 FIGURE 8d shows a flow chart exemplifying the iteration of the process as demonstrated in Figure 8c for any peptide P1 with a known sequence.

FIGURE 9 shows a table of the monoisotopic and the average mass changes of a peptide upon the respective modification.

20 FIGURE 10 shows a table with exemplified motifs of chemical and enzymatic reactions, their respective mechanism/enzyme and the resulting average mass difference of the modified peptide.

25 FIGURE 11 shows a table listing the most common amino acids, their three- and one letter codes as well as their monoisotopic and average mass in their dehydrated form.

FIGURE 12 shows a table listing the common amino-terminal and carboxy-terminal groups of peptides, as well as the chemical composition, their respective monoisotopic and average mass.

30 FIGURE 13 shows a table that refers to the fraction shifts of peptides that are caused by addition of the respective amino acid to the peptide sequence under the described experimental settings with cerebrospinal fluid as sample source.

35 FIGURE 14a shows a table with amino acids and their empirically derived occurrence before the N-terminal cleavage position (start position) of a peptide in a precursor sequence, the

respective overall occurrence of the given amino acid in all determined sequences and the ratio thereof.

FIGURE 14b shows a table with amino acids and their empirically derived occurrence after 5 the N-terminal cleavage position (start position) of a peptide in a precursor sequence, the respective overall occurrence of the given amino acid in all determined sequences and the ratio thereof.

FIGURE 14c shows a table with amino acids and their empirically derived occurrence before 10 the C-terminal cleavage position (end position) of a peptide in a precursor sequence, the respective overall occurrence of the given amino acid in all determined sequences and the ratio thereof.

FIGURE 14d shows a table with amino acids and their empirically derived occurrence after 15 the C-terminal cleavage position (end position) of a peptide in a precursor sequence, the respective overall occurrence of the given amino acid in all determined sequences and the ratio thereof.

FIGURE 15 shows a flow chart schematizing the procedural steps of an interaction of the 20 Differential Network Module with the CAN Module according to the present invention.

FIGURE 16 shows a flow chart schematizing the procedural steps of an interaction of the Marker Panel Network Module with the CAN Module according to the present invention.

25 FIGURE 17 shows a flow chart schematizing the procedural steps of an interaction of the Surrogate Network Module with the CAN Module according to the present invention.

FIGURES 18a and 18b show a table of the signal intensity values of the peptides with 30 coordinates Fraction 54; m/z 2743.0, Fraction 54; m/z 1371.5, Fraction 56; m/z 2927.2 and Fraction 20; m/z 1114.3 taken from 74 samples. Furthermore, the number of related peptides k with a Spearman's Rank Order Correlation Coefficient threshold of  $|r| \geq 0.8$  is shown.

35 FIGURE 19 shows a table with the measures of correlation of the signal intensities of the peptide with coordinates Fraction 54; m/z 2743.0 with some exemplary peptides using different measures of correlation.

FIGURE 20 shows a histogram of Spearman's Rank Order correlation coefficient probabilities. The value of a correlation coefficient of a peptide-to-peptide relation (x-axis) is plotted versus the probability for that peptide-pair to achieve that value (y-axis). Peptide-to-peptide pairs with low absolute correlation coefficients are most likely not related. This is 5 expressed by the maximum at zero of peptide-to-peptide relations from random data ( $P(r)$  Simulation). True positive relations are very likely to be found at higher absolute correlation coefficients. Therefore the plot of correlation coefficients of peptide-to-peptide relations from measured data ( $P(r)$  Measurement) deviates from  $P(r)$  Simulation, because correlation coefficients of functionally related peptides are most likely higher than those obtained from 10 random data. Such a plot should be generated when a threshold for a given CAN has to be chosen in order to exclude as much false positive peptide-to-peptide relations while including as many true ones as possible.

FIGURE 21 shows a table of identified peptides related to Chromogranin A 97-131, the 15 Spearman's Rank Order Correlation coefficient value of the said peptide with the related peptide, their relative monoisotopic mass and their amino acid sequence.

FIGURE 22 shows a graph exemplifying the usability of a Differential Network of the peptides SG I 88-132 and Chromogranin A 97-131. In hypothetic patients before 20 prostatectomy (black triangles) a correlation between these peptides is present ( $r = 0.97$ ), and a signal intensity ratio of about 10/1 is conserved. In hypothetic samples after prostatectomy (white squares) this ratio is not present and the Secretogranin I/Chromogranin A relation is "broken".

25 FIGURES 23a and 23b show a table of the signal intensity values of the peptides with coordinates Fraction 54; m/z 1371.5, Fraction 56; m/z 2927.2 and Fraction 20; m/z 1114.3 of 74 samples after removal of the variance of the signal intensity of the peptide with coordinates Fraction 54; m/z 2743.0. Furthermore, the number of related peptides  $k$  with a Spearman's Rank Order Correlation Coefficient threshold of  $|r| \geq 0.8$  after the removal of 30 said variance is shown.

FIGURE 24a shows a graph plotting the signal intensity of the peptide in fraction 54 with a 35 mass-to-charge-ratio of 2743.0 (F 54; m/z 2743.0) versus the signal intensity of the peptide in fraction 20 with a mass-to-charge-ratio of 1114.3 (F 20; m/z 1114.3). This plot exemplifies a pair of peptides showing no correlation.

FIGURE 24b shows a graph plotting the signal intensity of the peptide in fraction 54 with a mass-to-charge-ratio of 2743.0 (F 54; m/z 2743.0) versus the signal intensity of the peptide in the same fraction with a mass-to-charge-ratio of 1371.5 (F 54; m/z 1371.5). This plot exemplifies a correlation between a peptide-to-peptide pair consisting of a single charged and a double charged peptide ion.

FIGURE 24c shows a graph plotting the signal intensity of the peptide in fraction 54 with a mass-to-charge-ratio of 2743.0 (F 54; m/z 2743.0) versus the signal intensity of the peptide in fraction 56 with a mass-to-charge-ratio of 2927.2 (F 56; m/z 2927.2). This plot exemplifies a peptide-to-peptide pair exhibiting a functional relation.

FIGURE 25a shows a graph plotting the studentized signal intensity of the peptide in fraction 54 with a mass-to-charge-ratio of 2743.0 (F 54; m/z 2743.0) versus the studentized signal intensity of the peptide in fraction 20 with a mass-to-charge-ratio of 1114.3 (F 20; m/z 1114.3), i.e. the peptide pair of Figure 24a. A minimum spanning tree algorithm was performed to connect the nearest vertices. The path with the most vertices, i.e. the MST diameter, has been highlighted by a bold line. In this example, the path comprises 29 vertices.

FIGURE 25b shows a graph plotting the studentized signal intensity of the peptide in fraction 54 with a mass-to-charge-ratio of 2743.0 (F 54; m/z 2743.0) versus the studentized signal intensity of the peptide in the same fraction with a mass-to-charge-ratio of 1371.5 (F 54; m/z 1371.5), i.e. the peptide pair of Figure 24b. This plot exemplifies a correlation between a peptide-to-peptide pair consisting of a single charged and a double charged peptide ion. A minimum spanning tree algorithm was performed to connect the nearest vertices. The path with the most vertices, i.e. the MST diameter, has been highlighted by a bold line. In this example, the path comprises 50 vertices.

FIGURE 25c shows a graph plotting the studentized signal intensity of the peptide in fraction 54 with a mass-to-charge-ratio of 2743.0 (F 54; m/z 2743.0) versus the studentized signal intensity of the peptide in fraction 56 with a mass-to-charge-ratio of 2927.2 (F 56; m/z 2927.2), i.e. the peptide pair of Figure 24c. This plot exemplifies a peptide-to-peptide pair exhibiting a functional relation. A minimum spanning tree algorithm was performed to connect the nearest vertices. The path with the most vertices, i.e. the MST diameter, has been highlighted by a bold line. In this example, the path comprises 40 vertices.

35

FIGURE 26 shows a table where the peptides with coordinates Fraction 54; m/z 1371.5 and Fraction 56; m/z 2927.2 were tested according to a method of the present invention, whether

the given peptides are potentially n times charged ions of the peptide with coordinates Fraction 54; m/z 2743.0.

FIGURE 27 shows the precursor sequence of the "Hypothetical Precursor" (HP) using a 5 one-letter code. The sequence of the peptide HP 25-48 is underlined, the sequence of HP 25-50 is in bold letters.

FIGURE 28 shows an Averaged Peptide Display (peptide map) of CSF samples from 66 10 patients as discussed in Example 6. Each chromatographic fraction of each sample is analyzed by MALDI-ToF mass spectrometry; panel A: all 96 fractions generated from one sample are visualized as a "2-D gel-like picture" as shown in Figure 2. The x- and y-axis are 15 mass-to-charge ratio (m/z) and chromatographic fraction, respectively; grey-scale bars in the "2D gel-like picture" represent the peptide peaks where the intensity of the grey scale corresponds to the mass spectrometric signal intensity, which corresponds to the relative quantity of the peptide measured by MALDI. The insets B and C provide an enlarged view of the marked boxes in panel A, indicating that some peptides are present in more than one fraction.

FIGURE 29 shows the Correlation-Associated Network of VGF 26-58 (1) as the network hub 20 with VGF 177-191 (2), VGF 350-370 (3), VGF 26-59 (4, 5, adjacent fractions), VGF 23-59 (7), VGF 26-61 (8), VGF 26-62 (9), VGF 25-62 (10), VGF 485-522 (11) and VGF 373-417 (12-14, adjacent fractions). (6) is not a VGF peptide. Threshold of correlation was  $|r| \geq 0.68$ . The peptides are represented as bullets, and a peptide-to-peptide relation is displayed as a line, connecting two bullets. The network is projected onto the pertaining CSF peptide map 25 of Figure 28. The numbers in the brackets are also shown in the figure (without brackets).

FIGURE 30 shows the network members of Figure 29 mapped onto the protein precursor sequence of VGF. The numbers correspond to those of Figure 29. The arrows span the beginning and end of the respective sequences of the peptides. The VGF 26-58 Peptide 30 Network covers different parts spreading over the whole protein precursor.

FIGURE 31 shows the Correlation-Associated Network of Albumin 25-48 (1) as the network hub with Albumin 25-45 (2), Alpha-1-antitrypsin 397-418 (3), Albumin 25-48 (4,5, adjacent fractions), Albumin 27-50 (6,7, adjacent fractions) and Albumin 25-50 (8,9, adjacent fractions). Threshold of correlation was  $|r| \geq 0.67$ . The network is projected onto the pertaining CSF peptide map of Figure 28.

FIGURE 32 shows a table listing correct and false precursor predictions for any peptide-peptide relation as a function of threshold of correlation coefficient r.

FIGURE 33 shows the nomenclature used for peptide cleavages.

5

FIGURE 34 shows a table listing the numbers and percentages of amino acid residues found before amino-terminal, after amino-terminal (N+1), before carboxy-terminal (C-1) after carboxyterminal cleavage site (C+1) and at any position in all protein precursor sequences, e.g. Percentage(N-1)=n(N-1)/n(any position). Three out of the 139 peptides started at the beginning of the protein precursor, 34 peptides ended at the protein precursor sequence, therefore the sum of n(N-1) and n(C-1) deviate from 139. The percentage of an amino acid to be found in the respective positions is compared to the percentage at any position, indicating an x-fold increase or decrease of probability of a cleavage adjacent to that particular amino acid, e.g.  $x(N-1) = \text{Percentage}(N-1) / \text{Percentage}(\text{any position})$ .

10

FIGURE 35 shows a table listing the numbers and percentages of selected pairs of amino acids found before/after amino-terminal / carboxy-terminal cleavage sites and at any position in the protein precursor sequence. The percentage of a pair of amino acids to be found in the respective positions was compared to the percentage at any position, indicating an increase or decrease of probability, that this constellation influences cleavage.

15

FIGURE 36 shows a table listing two examples for predictions of peptide coordinates that were confirmed by the results of the ESI-MS/MS identifications.

20

FIGURE 37 illustrates the evaluation of the predictive power of the models: 139 peptides, all previously identified by ESI-MS/MS, were split into two groups. The first group of 70 peptides predicted the putative sequence of the second group. The peptide sequence information of the second group was suppressed during the prediction process. After completion of calculations, the putative proposals were confirmed by the results of the ESI-MS/MS identifications. The 139 different peptides corresponded to 224 different peptide coordinates on the Peptide Mass Fingerprint, since abundant peptides are present in more than one fraction.

25

FIGURE 38 shows a table listing percentages of correct proposals of precursor proteins and start-stop positions for six different models separately assessed for any stored proposal and for the proposal with most bonus points. At  $|r| \geq 0.75$ , proposals were generated for 27 out of

30

112 peptide coordinates. Since any peptide coordinate can store a list of up to three proposals, 81 proposals were generated.

FIGURE 39 shows a table with the peptides correlating to the peptide albumin 25-48 of CSF

5 samples taken from patients with different severe disruptions of the blood-CSF barrier.

Shown are the correlations  $r$  to the albumin quotient (complet albumin protein) of patients with damaged brain barrier (example 7, Figure 40), correlations  $r$  to the albumin 25-48 mass spectrometric signal intensity of patients with damaged brain barrier (example 6), the name of the peptides (Albumin 25-48, 27-50, 25-51 and alpha-1-antitrypsin 397-418) the

10 theoretical monoisotopic mass of the peptides and the sequences of the peptides.

FIGURE 40 shows five plots of the relative MALDI singnal intensities of the peptides Albumin

25-48, 27-50, 25-50, 25-51 and alpha-1-antitrypsin 397-418 relative to the albumin quotient

measured in patients with damaged brain barrier (example 7). In all cases there is a near

15 linear relation between the MALDI signal intensity and the albumin quotient.

#### Detailed Description of the Invention

Prior to giving a detailed albeit exemplary description of embodiments of the present

invention the following definitions are provided to establish how the technical terms are to be

20 understood herein.

#### **Definitions**

Unless defined otherwise, all technical and scientific terms used herein have the meaning commonly understood by a person skilled in the art to which this invention belongs. As used

25 herein, the following terms have the meaning ascribed to them unless specified otherwise.

"Sample" refers to any material, substance or the like containing or potentially containing peptides.

30 "Peptides" refers to polymers of amino acids coupled by peptide bonds comprising at least two amino acids. These amino acids can be the twenty standard amino acids and additionally unusual amino acids as known in the arts including D- and L- amino acids. Peptides can contain additional modifications such as posttranslational, enzymatic and/or chemical modifications.

35

"Status of a sample or an organism" means that the status or type of a sample at the time of the generation of the sample, e.g. the drawing of blood, is reflected by the contents and the

activities of the sample. The actual status of an organism at the time of sample generation (such as the drawing of blood) is reflected in the contents and activities present in the sample. The sample conserves the status similar to a snap shot picture. A status for example can represent the presence or absence of a certain disease, the presence or 5 absence of pregnancy, the sex of the individual from which the sample originated, the presence of a certain genetic variation such as the knockout of a gene or a polymorphism, the over expression or increased activity of a certain gene or gene product (for example as a consequence of a drug or of the transfection of the gene coding for the gene product or by direct addition of the gene product, etc.), the suppression of the expression or activity of a 10 certain gene or gene product (for example as a consequence of a drug, anti-sense nucleotides, RNAi (RNA interface) nucleotides, ribozymes, triplex-forming nucleotides, antibodies, etc.), the presence of genetic modified ingredients in food, cosmetics or other products, the age of the organism from which the sample originated, the species of the organism from which the sample originated, a certain treatment of the organism from which 15 the sample originated (for example with a therapeutically active substance a food ingredient or substance present in cosmetics, treatment with insecticides, pesticides or other toxic substances, etc.), the geographic origin of the sample, the development stage of the organism from which the sample originated (for example the stage of a fertilized egg, an embryo, an adult, intracellular/extracellular bacteria/virus, egg/larva/pupal/adult-stage of for 20 example butterflies, different development stages of plasmodium, etc.), the metabolic state of the organism from which the sample originated (for example hibernation, stages of the circadian rhythm, etc.), the point of time before, during or after the treatment of an organism with a substance, the localization (or tissue) within the organism, from where the sample was taken, and the like.

25 "Measurement parameter of a peptide" refers to any parameter known to or measurable by the investigator such as the molecular mass of the peptide, the mass/charge ratio of the peptide, the signal intensity of the measured peptide, the actual concentration of the measured peptide, the fraction-number in which the peptide is present as a consequence of 30 a certain separation protocol subjected to the sample, or the measured activity of the peptide.

35 "Correlation" or "relation" refers to a hypothesized mutual dependency of at least one parameter of two peptides, may this dependency be symmetric or asymmetric, known or not known, statistically significant or not. Relations of two peptides can be caused by chemical and biochemical reactions from one peptide to the other, by concerted gene regulation of the analytes, by common precursor peptides and so on.

"Measure of correlation", "correlation measure" or "measure of association" refers to statistical means to describe the symmetric or asymmetric statistical dependency of measurement parameters of pairs of peptides in terms of their "relation". Examples for measures of correlation are: "Pearson Product-Moment Correlation Coefficient", 5 "Spearman's Rank-Order Correlation Coefficient", "Kendall's Tau", "Kendall's Coefficient of Concordance", "Goodman and Kruskal's Gamma", "Manhattan distance", "Euclidean distance" and "Minimal Spanning Tree Diameter".

"Correlation associated network (CAN)" refers to the complete network of all measures of 10 correlation identified within samples representing one status or identified within different groups of samples representing different statuses. It is possible that more than two peptides correlate to each other and a CAN contains at least two peptides correlating to each other. It should be noted that the peptide CAN based on "a sample" does not necessarily comprise results obtained from a single experiment. Rather, to completely determine a peptide CAN, 15 multiple experiments are often needed, and the combined results of which are used to construct the peptide CAN for that particular sample. The results of the calculation of a CAN (CAN of first order) can be used for another round of calculation of measures of correlation and so on. The results of these kind of calculations are also termed CANs or more specifically CANs of second or higher order.

20 "Peptide-topology" refers to the entirety of measured and computed peptide data of a sample ("measurement parameter of peptides") comprising the masses of the peptides, the signal intensity of the peptides (preferably measured by mass spectrometry or another measurement method suitable to quantify peptides), the fraction number (if the sample was fractionated prior to mass spectrometry) and measures of correlation calculated using these data.

"Groups of samples" refers to a set of samples corresponding to a certain status. A group of 30 samples for example could comprise 10 plasma samples of diabetic patients. The samples of a group need not to be of exactly the same origin. For example a group of samples may also comprise 5 plasma samples of diabetic patients and 5 urine samples of diabetic patients. The reason for this being that many peptides present in plasma are also present in urine and for example the same diabetes-specific peptides may be present in plasma and urine, as long as the sample originates from a diabetic patient.

35 "Known peptide" means that the peptide with that particular sequence or part of a sequence in the sample is known to the user of the invention. An unknown peptide is a peptide whose

sequence is not known to the user of the invention, although the sequence of the peptide may be known from the literature or other sources of information such as sequence data bases.

5 "Potential peptide" refers to a mass spectrometric signal which most likely represents a peptide.

"Precursor of a peptide" refers to the longest amino acid sequence present in nature comprising the sequence of the peptide, i.e. form which the peptide can originate.

10 "Coordinate(s) of a peptide" refer to the mass-to-charge ratio and optionally further specific measurable properties obtainable by a detection or identification process that are involved in the identification and/or quantification of the said peptide/peptide ion. In the examples of this invention the peptide coordinates are the elution time/fraction number of a chromatographic 15 process and the mass-to-charge ratio, thus comprising of two coordinates. In this invention, these coordinates often are written in a short form, such as "F 56; m/z 2873.0", which identifies the signal of a peptide found in fraction 56 with the mass-to-charge ratio 2873.0. Of course, further dimensions can be necessary, such as a previous capillary electrophoresis, or a downstream second mass spectrometric process. "Coordinates of a peptide", "signal 20 coordinates" or "peptide" are often used synonymously.

"Fitness value" refers to an assessment of a predicted sequence based on the experimental properties of an unknown peptide. Any predicted sequence gains points for properties that match the experimental properties, such as the correct prediction of the fraction number. 25 The higher the "fitness value", the more probable the correctness of the predicted sequence. According to the present invention fitness values are manually or automatically suggested for each sample type and empirically tested for suitability.

"Landmark peptides" refers to peptides that are related to numerous other peptide signals 30 and least related to each other. The identification, e.g. sequencing, of these landmark peptides should be prioritized to gain a rapid overview about the peptide composition of a sample.

#### Providing the data

35 Figure 1 shows schematically the hardware components and software modules according to the present invention, their interfaces and the flow of information between the hardware components and the software modules. Although measurement data could be provided

without fractionating the samples prior to performing mass spectrometry, such a fractionation e.g. by chromatography into e.g. 96 fractions is preferred. A "fraction" in terms of chromatography is a part of the effluent recovered during a separation step. Usually, several fractions are collected. Fractions usually contain different "subsets" of peptides from the

5 sample. Suitable separations methods for peptides are chromatographic methods such as ion exchange, hydrophobic interaction, iso-electric focusing, gel filtration or affinity chromatography, electrophoretic methods such as native, iso-electric, denaturizing or SDS-gel electrophoresis using matrixes such as polyacrylamide or agarose gels, paper electrophoresis, thin layer chromatography, capillary electrophoresis, methods using

10 centrifugation for separation such as sucrose or Caesiumchloride gradient centrifugation and the like. These chromatographic fractions are then subjected to a measurement of the mass spectrum providing 96 mass spectra, which can be visualized e.g. in a 2D gel-like format as shown in Figure 2. To this end all kinds of methods suitable to determine masses of peptides and preferably all kinds of mass spectrometry can be used in the present invention such as

15 matrix-assisted laser desorption time of flight (MALDI-TOF) mass spectrometry, liquid chromatography electro-spray ionization (ESI) quadrupole time of flight mass spectrometry (LC-ESI qTOF) and the like. It is furthermore possible to analyze only selected but not all fractions by mass spectrometry.

20 Each bar in Figure 2 depicts a peak present in one of the 96 mass spectra, wherein the colour intensity of the bar corresponds to the intensity of the corresponding mass spectrometric signal. The x-axis of Figure 2 represents the mass to charge ratio  $m/z$  and the y-axis the chromatographic fraction number. The  $m/z$  values preferably range from 1:000 to 15.000, although higher or lower  $m/z$  values can be included as long as these values can be

25 resolved by mass spectrometry or other methods. Within this range of  $m/z$  values the detected peptides can be comprised of only two amino acids at the lower end up to peptides of very large molecular mass such as alpha-2 macroglobuline having a molecular mass of 725 kDa.

30 Similar 2D gel-like maps are produced for every sample out of the set of samples to be analyzed. These maps can be averaged yielding an averaged peptide mass fingerprint map as shown in Figure 2. This averaged map serves as a template for the definition of usually about a thousand peak coordinates, i.e. the x-coordinate corresponding to the  $m/z$  value and the y-coordinate corresponding to the fraction number. In practice one selects those peak

35 coordinates exhibiting a signal above a certain threshold value.

### Data Pre-processing

In order to obtain measurement data that is suitable for a correlation analysis and that gives meaningful results preferably a pre-processing of the data is performed using methods such as baseline correction, spectra normalization, outlier detection and the like. Methods for 5 baseline correction are well known in the art (e.g. Fuller et al, *Applied Spectroscopy*, 42, 217 1988). In a preferred embodiment the pre-processing of the data is performed by applying the baseline correction being part of the software RAZOR Library 4.0, Spectrum Square Associates, Ithaca NY, USA. Optionally a normalization of the mass spectra can be performed by using the signal intensities or the integrated mass spectra. Outlier samples can 10 be identified by means of a principal component analysis as provided by the commercially available software package Pirouette 3.0, Infometrix Inc., WA, USA. Based on this principal component analysis individual mass spectra or even whole samples that exhibit a Mahalanobis distance  $M_D$  above a critical threshold value of should not be considered for a further analysis and thus be discarded. In the examples described further below a 15 Mahalanobis distance of  $M_D > 11.5$  was chosen for 74 samples.

The preprocessing, processing and display of the data according to the present invention can be performed e.g. on a Apple G4 Computer, wherein the CPU consists of 2 processors with 800 MHz each and the memory size is 1.25 Gigabyte. The local data storage of peptide-20 to-peptide relations (measures of correlation, coordinates of peptides) is accomplished by a local Valentina Database system (Valentina 1.9 for Realbasic, Paradigma Software, Beaverton, Oregon, USA). The peptide sequence information is provided by a proprietary Interbase Server (Interbase 6, Borland Software Corp., Scotts Valley, CA, USA). Microsoft Internet Explorer 5.1 for Apple computer systems can be used for representation of results 25 from internet resources. The CAN software launches the Internet explorer with a specific address that contained the keywords for querying the Swiss-Prot, the PubMed, and the US Patent database. Visualization of three-dimensional objects can be performed using a Realbasic RB3D engine (RealBasic 3.5, RealSoft, Austin, Texas, USA).

30 Other digital computer system configurations can also be employed to perform the methods of the present invention, and to the extent that a particular system configuration is capable of performing the method of this invention, it is equivalent to the representative digital computer system schematically shown in Figure 2. Once they are programmed to perform particular functions pursuant to instructions from program software that implements the methods of the 35 present invention, such digital computer systems in effect become special-purpose computers particular to the methods of this invention.

Computer programs implementing the methods according to the present invention will commonly be distributed to users on a distribution medium such as floppy disk or CD-ROM. From there, they will often be copied to a hard disk or a similar intermediate storage medium. When the programs are to be run, they will be loaded either from their distribution 5 medium or their intermediate storage medium into the execution memory of the computer, configuring the computer to act in accordance with the method of this invention. All these operations are well-known to those skilled in the art of computer systems. The term "computer-readable medium" encompasses distribution media, intermediate storage media, execution memory of a computer, and any other medium or device capable of storing a 10 computer program implementing the methods of this invention for later access by a computer.

#### **Correlation Associated Network Module**

As exemplified by the arrows in Figure 1 the raw measurement data or preferably the 15 preprocessed measurement data is supplied to the so called Correlation Associated Network (CAN) Module 42. Of the modules 40 of the present invention the CAN Module 42 is the most fundamental one. Basically the CAN Module 42 scans the measurement data obtained for example from Liquid Chromatography-Mass Spectrometry (LC-MS) experiments 22. On the basis of this data correlations of peptides are searched for by calculating measures of 20 correlation between for example their relative concentration as measured by mass spectrometry.

Measures of correlation can be used to represent the degree of relationship between two variables throughout many observations. These variables can be either correlated, not 25 correlated or anti-correlated. In the context of the present invention measures of correlation are used to detect such correlated, not correlated or anti-correlated peptides in a set of samples. This can be done e.g. by calculating Spearman's rank-order correlation coefficient of the signal intensities of two peptides measured in several samples. Preferably this is done for all pairs of peptides. Once these measures of correlation have been calculated only those 30 pairs of peptides are selected that exhibit a certain behaviour, i.e. a certain degree of correlation, a certain degree of anti-correlation or a certain degree of no correlation at all. The parameters of such selected peptide pairs, e.g. the coordinates of the two peptides of each peptide pair, the measure of correlation, etc., can be stored, displayed on a display device or further processed. Preferably the data is stored in a database, as a text file or in 35 another computer-readable form. Alternative measures of correlation to Spearman's rank order correlation coefficient are Pearson Product-Moment Correlation Coefficient, Kendall's

Tau, Kendall's Coefficient of Concordance, Goodman and Kruskal's Gamma and Minimal Spanning Tree diameters.

A Minimal Spanning Tree (MST), also known as Minimum Spanning Tree, is defined by the 5 collection of edges that joins together all points in a connected set of data points, with the minimum possible sum of edge values (e.g. Evan, *Graph Algorithms*, Computer Science Press, 1979). An edge can be graphically displayed by a line connecting two data points. A MST can be graphically displayed by a set of points (data points) connected by a minimum of lines to each other. Examples of MSTs are shown in Figures 25a to 25c as described in 10 more detail further below. A MST also provides a plausible "connectionist" approach to solving the "Traveling Salesman" problem (e.g. Kruskal, *Proc. American Math. Soc.*, 7, 48-50, 1956; Sun et al, *Physica A*, 199, 232-242, 1993), which identifies the minimum connected path between all data points. The MST diameter can be defined as the maximum 15 number of edges in the paths of a graph of a MST. Usually a correlation, for example a Spearman's rank order correlation coefficient, is used to find a measure of correlation or association or dependency between variables, i.e. data points. A problem is that correlation is sensitive to linear trend, and linear trends are not always well presented for two associated variables. In the present invention, the diameter of the MST is used as an alternative 20 measure of correlation between two variables. In order to use the diameter of the MST to analyze a given set of n statistical observations, all observations should be connected via the MST and then the MST diameter should be calculated. The larger the MST diameter is, the stronger is the association between two variables. In the context of mass spectrometry signal intensity data (in the present invention preferably MALDI mass spectrometry signal intensity data) it was found, that MST diameters  $> 0.425$  times n indicate a noticeable 25 association between signal intensities of peptide coordinates. In general all kinds of mass spectrometry signal intensity data, such as MALDI or ESI mass spectrometry data, can be used according to the present invention.

As already mentioned pairs of peptides are tested for their degree of correlation by 30 estimating e.g. Spearman's rank order correlation coefficients between their signal intensities throughout many observations. It turns out that pairs of peptides which are biologically or functionally related surprisingly often exhibit correlation coefficients that are much higher than correlation coefficients that would be expected by chance. Unrelated pairs of peptides have low absolute values of correlation coefficients. Figure 3 exemplifies the 35 correlational behaviour of related peptides. Four traces of spectra from four different samples are focused on the signals of a human osteopontin peptide being comprised of the amino acids 249-314 of human osteopontin ( $m/z = 7653.6$  Da) and its phosphorylated

derivatives, carrying one ( $m/z = 7733.5$  Da), two ( $m/z = 7813.5$  Da), three ( $m/z = 7893.4$  Da) or more phosphorylated residues. The conserved concentration ratios of the peptides between samples leads to high degrees of correlation of the signal intensities of the respective peptide pairs.

5

Using the results of the above described computations of measures of correlation so called correlation associated networks (CANs) can be defined. A CAN, i.e. a network of peptide relations, comprises a peptide of interest, the so called hub peptide, and all those peptides and sample parameters that correlate to a certain degree with the hub peptide. The term hub 10 is used in a similar manner in the theory of network topology and is to characterize the resemblance of a hub peptide to the hub of a wheel, the hub peptide being at the center of the spokes representing the peptide-to-peptide relations and the correlating peptides being at the respective ends of the spokes. In practice, the composition of a CAN is highly dependent on the threshold of correlation as selected by the user. This threshold is chosen 15 according to the goal of a user. If a user is searching for peptides that strongly relate to a peptide of interest, such as peptides stemming from the same precursor, then he will select a threshold that will cause a selection of only the upper 5 % of the strongest correlations with the peptide of interest. The threshold value to be chosen for e.g. the Spearman's rank order correlation coefficient depends for the thus selected subset on the number of samples and 20 the peptide of interest. In case the user is interested in finding functionally related peptides, such as e.g. peptides being co-secreted from vesicles, the user will choose a threshold value, that will select e.g. the upper 10 % of the strongest correlations.

The hub peptide and the peptides related thereto and selected as described above represent 25 a CAN of first order. Depending on the objective it can be necessary to compute CANs of higher order due to the complexity of biological networks and pathways. As explained above, CANs connect directly related peptides which exhibit a high degree of correlation. Adjusting the threshold to lower values results in including more loosely related peptides into the network as well as increasing the probability of predicting false relations. For this reason a 30 preferred embodiment of the present invention contemplates the computation of CANs of higher orders, such as e.g. second and third order. Since the direct members of a network of interest constitute the first order neighborhood, all these members are potential starting points for the calculation of second order neighborhoods as shown in Figure 4. Although computing CANs of higher order certainly will improve the results, the computational 35 requirements set an upper limit, because the computational effort increases with the order of the CANs. A calculation of a CAN of the  $n^{\text{th}}$  order, where  $n$  is greater than 5, can require more than several millions of calculations. Thus, this approach preferably should be used for

the analysis of rather complex samples in order to include indirectly related peptides, thus avoiding having to decrease the value of the correlation threshold and possibly including false relations.

5 For any kind of sample the composition of peptides varies, novel peptide coordinates emerge, others disappear and many peptide coordinates have a different peptide sequence aligned to it. This results in dealing with many unknown peptide coordinates when operating with novel sample sources (types of samples). In order to accelerate analysis of an interest list or more general, to analyze the overall peptide composition of a sample, according to the  
10 present invention it is possible using CANs to accelerate the identification of peptides in complex biological samples by defining a list of representative peptides, so called landmark peptides, for further analysis such as peptide sequencing based on CANs described further below. The method comprises the following steps, which are shown in Figure 5. At step 80 mass spectra are provided as described above, wherein the peaks of the signal intensities in  
15 the mass spectra correspond to potential peptides. Then the measures of correlation between the measured signal intensities corresponding to potential peptides are computed (step 82). Thereafter at step 84 those peptides are grouped together that exhibit a degree of correlation above an adjustable threshold. These selected peptides constitute a CAN present in the samples analyzed. Finally, one peptide out of each determined CAN is assigned to  
20 represent that respective CAN at step 86. In doing so a plurality of landmark peptides is provided being representative of the analyzed samples. These landmark peptides have the properties of being hub peptides and they are least related to other peptides within the same type of samples. Identifying a list of these landmark or prioritized peptides gives a rapid overview about the peptide composition present in complex biological samples, omitting the  
25 majority of similar peptides from for example the same precursor peptide. This is useful to obtain a general overview of the key peptides present in a sample or present in an interest list from that sample.

It is contemplated that a such generated interest list of prioritized landmark peptides will  
30 contain a set of n peptide coordinates, and for any peptide z the number of relations, which the peptide z has at a defined threshold r,  $k_{z,r}$ , will be determined. The peptide z with the highest value of  $k_{z,r}$  will be defined as y and be rank 1 etc. on the prioritization list. Then the variance of signal intensities of that determined peptide coordinate y will be removed from the signal intensities of the related peptides x in a data matrix, for example by a combination  
35 of formulas 1, 2 and 3 shown below. Then this peptide will be removed from the data matrix. Calculations of any k and r start from the beginning to determine the representative peptide ranked number 2 in the prioritization list, and so on. Calculations end for example when the

data matrix contains no more peptide coordinates, or no peptide has more than zero relations, or the number of peptide coordinates desired has been reached.

**Formulae 1 to 3: Removal of variance of the representative peptide coordinate y on**

5 **peptide coordinate x**

$$X_{VR,p} = X_p - a_{xy} - b_{xy} Y_p \quad (1)$$

where

$X_{VR,p}$  : Signal intensity of peptide x at observation p, Variance of peptide y removed

$X_p$  : Signal intensity of peptide x at observation p

10  $Y_p$  : Signal intensity of peptide y at observation p

$m$  : number of observations

$$a_{xy} = \frac{\left( \sum_{p=1}^m X_p \right) \left( \sum_{p=1}^m Y_p^2 \right) - \left( \sum_{p=1}^m Y_p \right) \left( \sum_{p=1}^m X_p Y_p \right)}{m \sum_{p=1}^m Y_p^2 - \left( \sum_{p=1}^m Y_p \right)^2} \quad (2)$$

15 and

$$b_{xy} = \frac{\sum_{p=1}^m X_p Y_p - \frac{1}{m} \left( \sum_{p=1}^m X_p \right) \left( \sum_{p=1}^m Y_p \right)}{\sum_{p=1}^m Y_p^2 - \frac{1}{m} \left( \sum_{p=1}^m Y_p \right)^2} \quad (3)$$

It is further contemplated that peptides being part of a CAN preferably are represented by

20 graphical objects such as e.g. bullets and their mutual relations by lines connecting these bullets. In order to enable a more intuitive analysis of the results, this network can be projected onto a peptide map as shown in Figure 6. Peptides that have been identified can be provided with links to databases, with lists containing additional information about these peptides or with other sources of additional information regarding said peptides.

25

The coordinates or measurement parameters of related peptides can be queried in public, commercial and/or proprietary databases in order to identify further data about the potential identity, function or use of the corresponding peptides. Suitable public databases include e.g. the PubMed literature database, the OMIM disease database, the NCBI-Sequence

database (all provided by the US National Library of Medicine, MD, USA), the Swiss-Prot and TrEMBL Sequence database, enzyme database, Swiss 3D image database, Prosite protein family and domain database (all provided the Swiss Institute of Bioinformatics, Switzerland), patent databases of the US, European, Japanese, German patent offices, the 5 Gene Cards database of the Weizmann Institute, etc. Suitable commercial databases are for instance commercial patent databases containing patented amino acid or nucleic acid sequences such as DGENE (Thomson Derwent, USA) or REGISTRY (Chemical Abstracts Service, USA). A suitable proprietary database is the database of the user containing peptide sequences from various sources and species. This combination of the visualization 10 of peptide networks and the connection to many sources of information alleviates the evaluation of the identified peptides for potential uses such as their use as therapeutic peptides or as biomarkers as will be described in more detail further below.

As is apparent from the above, correlation associated networks can be used to generate 15 hypotheses about relations between structurally and/or biologically related peptides. These hypotheses are based on a correlational analysis of signal intensities and corresponding relative peptide concentrations from independent samples. The examples described in the sections further below will demonstrate that correlation associated networks are powerful 20 tools for the systematic analysis and interpretation of large peptidomic and proteomic data in order to reveal functional relationships governing protein synthesis, posttranslational modifications and degradation. CANs support the discovery of novel bioactive and diagnostic peptides leading beyond the mere comparison of peptide concentration changes caused by a disease.

25 According to the present invention the CAN Module 42 is interacting with several application modules 44 comprising a Sequence Network Module 46, a Differential Network Module 48, a Marker Panel Network Module 50 and a Surrogate Network Module 52 as shown in Figure 1. These application modules 44 of the present invention and their interaction with the fundamental CAN Module 42 will be described in detail in the sections below.

30

#### **Sequence Network Module**

The interaction of the Sequence Network Module with the fundamental CAN Module according to the present invention allows to predict the amino acid sequences of unknown peptides with or without modifications of the sequence and/or to predict unknown 35 modifications of a known or unknown peptide sequence. Although the identity of the peptide is unknown, certain physicochemical and biochemical properties of the signal of an unknown peptide are known and can be exploited for amino acid sequence prediction such as the

mass-to-charge ratio (m/z) or the chromatographic behaviour (fraction number/retention time). Furthermore bioinformatic support data shown at 56 in Figure 1 such as the correlation associated network of related peptides, mass differences and differences in fraction number between the peptides of the correlation associated network, and the like are 5 accessible as they can be computed using experimental data and the amino acid sequences of other members of the correlation associated network possibly already known.

Figure 7 shows a flow chart schematizing the procedural steps of an interaction of the Sequence Network Module with the CAN Module according to the present invention allowing 10 the prediction of the sequence of peptides using the peptide topology of a plurality of samples containing a peptide having a known precursor. At step 80 a respective mass spectrum for each sample of said plurality of samples is provided, wherein the signal intensity peaks correspond to potential peptides. Thereafter at step 88 the peptide having a known precursor is identified using the mass of said peptide, wherein the sequence of the 15 known precursor is known. Then measures of correlation between the signal intensity of the peptide having a known precursor and the signal intensities of the other potential peptides are computed at step 90. At step 92 those potential peptides are selected, which exhibit a degree of correlation with the peptide having a known precursor above a certain adjustable threshold, and finally the sequence of the potential peptides are predicted at step 94 by 20 matching masses of putative fragments of the sequence of the known precursor with the masses of the potential peptides correlating with said peptide having a known precursor.

Alternatively after step 92 the mass differences between each of the potential peptides and the known peptide can be computed at step 96, and thereafter the sequence and/or the 25 biologically, chemically or physically modified sequence of the potential peptides predicted at step 98 by using data about mass differences caused by biological, chemical or physical processes matching the mass differences determined in step 96.

The first of the above approaches is more comprehensive, since all plausible putative 30 sequences are generated from the precursor sequence of the known peptide (steps 90-98). The second approach (steps 90-96, 100-102) generates fewer but more reliable predictions. It has been observed that related peptides very often have very similar sequences/sequence 35 modifications, and these predictions are promoted by the second approach. Nevertheless, since both approaches have steps 90-96 in common, computational power is "saved" if both approaches are combined in one operation, as contemplated in the present invention.

Mass differences may result from addition or removal of N- or C-terminal amino acid residues or of posttranslational modifications of amino acid side chains such as phosphorylation, amidation, sulfatation, glycosylation, fatty acids or Ubiquitin modification, and the like or chemical modifications such as oxidation, disulfide bonding, and the like or N- or C-terminal modifications such as pyroglutamate modifications and the like. All of these modifications result in distinct increases or decreases of the molecular mass of the corresponding peptide. Also internal insertions or deletions or the exchange of one amino acid for another, so called point mutations, result in exactly predictable mass changes of the peptide.

10

According to the present invention the prediction of sequences is possible regardless of whether the identity of one of the related peptides is known or not. Especially if the identity of one peptide is known, mass differences corresponding to the molecular masses of amino acid residues allow to directly predict the complete sequence of the unknown peptide with high reliability. If the identity of no peptide is known, than for example it can be predicted that the unknown peptide 1 and the unknown peptide 2 are identical, except that for example peptide 2 contains an additional amino acid residue, for example a Tyrosine residue, or for example peptide 2 is the same peptide as peptide 1 except that it is phosphorylated, etc. The prediction is not always correct, but the more independent information is accessible, the more reliable the prediction becomes. For example if the mass difference fits to the addition of an Tyrosine amino acid residue and the peptide is present in a fraction, which fits to the prediction of the fraction-shift of a peptide with an additional Tyrosine residue, the overall reliability of the prediction increases.

25 For this embodiment the use of proprietary and/or commercial and/or public databases is possible. Suitable databases are for example databases containing amino acid or nucleic acid sequence information such as the NCBI sequence data base, Swiss-Prot, the EMBEL sequence data base, the DNA data base of Japan, data bases of patented sequences, and the like, data bases containing information about the structure of carbohydrates, such as 30 PROSITE (Falquet et al, *Nucleic Acids Res.*, 30, 235-238, 2002), data bases containing information about posttranslational, enzymatic or chemical peptide modifications such as phosphorylation sites of peptides, glycosylation sites of peptides, positions of unusual amino acids such as hydroxy-proline or hydroxy-lysine within peptides, databases containing information about recognitions sites of proteases, ligases, phosphatases, kinases, and the like within peptide sequences, databases containing information about the susceptibility of 35 certain amino acids or sequences of amino acids towards chemical modifications such as oxidation, reduction, intra-molecular rearrangement, data bases containing data about three-

dimensional structures about peptides, carbohydrates or other biological structures, and the like (Falquet et al, *Nucleic Acids Res.*, 30, 235-238, 2002). All of these different kinds of databases enable to predict the structural difference between peptides, based on certain incremental increased or decreased molecular masses of these peptides. For example:

5

(i) amino acid sequences stored in databases allow the calculation of the masses of successive shortened or extended peptides or of peptides containing mutations of their sequence

10

(ii) databases containing for example recognition sites (sequences) of kinases allow to predict, that the molecular mass of a certain peptide, containing such a recognition site, may have a molecular weight increased or decreased by the mass of a phosphate group

15

(iii) data bases of recognition sites of proteases allow to predict the molecular masses of potential proteolytic fragments of a certain peptide

20

(iv) databases containing experimental data about physical properties of peptides such as elution times during for example hydrophobic interaction chromatography allow to predict, if a certain peptide sequence with a certain molecular mass is likely to elute at a certain time point during chromatography

25

(v) databases containing prediction values of chromatographic retention times or fraction numbers based on the amino acid composition and/or sequence of the peptide: if a certain chromatographic column is used, a peptide with an additional tyrosine residue would elute 3 fractions later than a peptide without that additional tyrosine residue. For example a peptide I in fraction x with mass y is known and a related peptide II within fraction x+3 has the molecular mass y plus the mass of a tyrosine residue. This would indicate with high reliability that peptide II is the same peptide as peptide I, except that it contains an additional tyrosine residue somewhere within its sequence

30

(vi) databases of three-dimensional structures of for example peptides allow to predict, if there is for example space enough at a certain amino acid side chain to be modified for example by a phosphate group or a sugar moiety, resulting in an increased molecular weight of the potential corresponding peptide

35

The prediction of physicochemical and biochemical properties of putative amino acid sequences fit surprisingly well to experimentally determined properties. This approach can be extended by utilizing knowledge about precursor amino acid sequences and posttranslational, chemical and enzymatic modifications of known related peptides as 5 provided by the support data 56 shown in Figure 1 and as discussed above. Furthermore, information about a known peptide such as the name of its precursor, its precursor sequence, its start and end-position within the precursor sequence can be retrieved before or during the prediction processes. Information about protease recognition sites, predictions 10 of domains, and structures sensitive to proteolytic digestions can also be retrieved. This information can be supplied manually, from databases or lists or from a comparable source of information. A conversion of mono-isotopic m/z ratio to average m/z ratio, from the m/z ratio of the charged ion to the m/z ratio of the un-charged ion within a reasonable error tolerance is known to those skilled in the art.

15 The invention comprises specific rules, which determine if a putative amino acid sequence derived according to one of the methods described above fits to the peptide signal coordinates of an unknown peptide. These rules which are schematically shown in Figures 8a to 8d can be applied in any order and it is not necessary to apply all of them in any given case:

20

Rule a:

This rule applies formula 4 (shown below) to check, whether the unknown peptide coordinate is an n-fold charged ion of the known peptide coordinate by the following condition, where n can be an integer number greater than 1,  $m/z_{\text{unknownpeptide}}$  is the m/z ratio of the unknown 25 peptide,  $m/z_{\text{knownpeptide}}$  is the m/z ratio of the known peptide and  $\text{Mass}_{\text{threshold}}$  is a maximum difference of the calculated mass from the measured mass. A preferable  $\text{Mass}_{\text{threshold}}$  equals the mass precision of the instrument and the subsequent data processing routines. If this condition is met, the proposal is rewarded with a high fitness value and the proposal that the unknown peptide is the n-fold charged ion of the known peptide can be stored.

30

**Formula 4: Check for n times charged peptide ions**

$$|n * (m/z_{\text{unknownpeptide}} - 1) - m/z_{\text{knownpeptide}}| = \text{Mass}_{\text{deviation}} \leq \text{Mass}_{\text{threshold}}$$

wherein the asterix (\*) indicates the mathematical operation of multiplication.

35

Rule b:

If the difference of the masses of a known hub peptide P1 and a related peptide P2 corresponds to a mass of an post-translational modification, as listed for example in the table "Mass Changes Due to Post-translational Modifications of Peptides and Proteins" shown in

5 Figure 9 or as known from the prior art (Falquet et al, *Nucleic Acids Res.*, 30, 235-238, 2002), then P2 is proposed to be the post-translationally modified derivative of P1. If the amino acid sequence of the known hub peptide P1 contains specific sites for posttranslational modifications or it is known that P1 is or can be posttranslationally modified, and if the mass difference between the known and the unknown peptide corresponds to the  
10 mass difference resulting from the presence or absence of that posttranslational modification, the fitness value is increased. The table shown in Figure 15 exemplifies motifs, enzymes recognizing these motifs and the resulting mass differences. Numerous other posttranslational modifications or putative sequence motifs bearing certain post-translational modifications are known in the prior art and could be used as well such as N-glycosylation or  
15 O-glycosylation sites (motifs), phosphorylation sites, sulfatation sites, and the like (e.g. Alberts et al, *Molecular Biology of the Cell*, Garland Publications, 2002; Coligan et al, *Short Protocols in Protein Science*, John Wiley & Sons, 2003; Falquet et al, *Nucleic Acids Res.*, 30, 235-238, 2002).

20 Rule c:

Putative sequences or putative fragments are generated from potential amino- and carboxy-terminal truncations or additions of amino acids of the known precursor sequence of the hub peptide and are checked whether they match the found m/z ratio of the unknown peptide coordinate. A putative sequence is generated by systematically and iteratively defining start-  
25 and end-positions, i and j, in the given precursor sequence of the hub peptide, as exemplified in Figure 8b. The mass of the putative amino acid sequence  $M_{CALC}$  is calculated by summing up the masses of the amino acids, the hypothesized posttranslational modifications of the amino acid residues and/or of the terminal groups of the putative amino acid sequence (see the tables in Figures 9, 11 and 11 and formula 5 shown below). Rule c defines that if the calculated mass differs from the measured mass  $M_{FOUND}$  of the unknown peptide signal by less than a given mass threshold  $T_{Mass}$ , this putative amino acid sequence plus posttranslational modifications are proposed and further rules d to i can be applied, otherwise this proposal is rejected. This can be done with one or more putative peptide sequences or with all hypothetically possible peptide sequences that can be deduced from  
30 the precursor sequence of the known related peptide signal coordinates.  
35

**Formula 5: Calculation of Masses**

$$M_{CALC} = n_A * M_A + n_R * M_R + n_N * M_N + n_D * M_D + n_C * M_C + n_E * M_E + n_Q * M_Q + n_G * M_G + n_H * M_H + n_I * M_I + n_L * M_L + n_K * M_K + n_M * M_M + n_F * M_F + n_P * M_P + n_S * M_S + n_T * M_T + n_W * M_W + n_Y * M_Y + n_V * M_V + M_{N-Terminal Group} + M_{C-Terminal Group} + M_{Modifications}$$

5

wherein:

$M_{CALC}$  is the calculated mass of the peptide with the given/putative sequence,

$M_{One\ Letter\ Amino\ Acid\ Code}$  is the mass of the appropriate amino acid,

$n_{One\ Letter\ Amino\ Acid\ Code}$  is the number of the appropriate amino acid in the given/putative 10 sequence,

$M_{N-Terminal\ Group}$  is the mass of the N-terminal group,

$M_{C-Terminal\ Group}$  is the mass of the C-terminal group, and

$M_{Modifications}$  is the mass change by modification(s), in the case of no modification

$M_{Modifications} = 0$ .

15

**Rule d:**

The number and the identity of amino acids influence the elution time/fraction number, depending on the size and the kind of the chromatography column used and the chromatography conditions. The fraction number/elution time of a peptide can be surprisingly 20 well predicted on the basis of its amino acid sequence by the so called Group Method of Data (GMDH, e.g. Mueller and Lemke, *Self-Organising Data Mining Extracting Knowledge From Data*, Trafford Publishing, 2003), multiple regression or comparable mathematic methods with a training set of peptides with known sequences, which are separated under the same chromatographic conditions as exemplified in Formula 6 shown below. In the said 25 training set, the number of any amino acid residue type of a peptide is the independent variable whereas the fraction number of the peptide is the dependent variable. If the calculated fraction number (e.g. Formula 6) of the predicted amino acid sequence matches the derived fraction number of the unknown peptide within a given error tolerance, then the model fitness points are increased. If the mass differences are proposed to be resulting from 30 distinct amino acid deletions/additions and if the differences in fraction number can be matched with these said amino acid sequence differences (see Figure 13), the model fitness points are increased.

**Formula 6: Estimation of fraction number based on proposed sequence**

$$F_{CALC} = 35.89 - 0.45 * n_S + 0.47 * n_E + 2.86 * n_I - 3.82 * n_H + 5.15 * n_L + 5.54 * n_F + 2.92 * n_Y - 1.72 * n_K - 0.85 * n_Q + 5.35 * n_W + 2.20 * n_V$$

wherein:

$F_{CALC}$  is the calculated Fraction number of the given sequence, and

$n_{ONE\ LETTER\ AMINO\ ACID\ CODE}$  is the number of the appropriate amino acid in the given sequence.

5 Rule e:

If the N-terminal position of the predicted amino acid sequence is the same as the N-terminal position of the known peptide, the fitness value is increased. This is because the known peptide and the unknown peptide of the underlying signals are related via a C-terminal proteolytic reaction, which is observed surprisingly often.

10

Rule f:

If the C-terminal position of the predicted amino acid sequence is the same as the C-terminal position of the known peptide signal, the fitness value is increased. This is because the known peptide and the unknown peptide of the underlying signals are related via an N-terminal proteolytic reaction, which is observed surprisingly often.

15

Rule g:

If the start position and/or the end-position of the predicted sequence is preceded or followed by sites of infrequent proteolytic events, the fitness value of this proposal is decreased. If the 20 start position and/or the end-position of the predicted sequence is preceded or followed by sites of frequent proteolytic events, the fitness value of this proposal is increased. This is because it has been observed that peptides are often products of specific and/or unspecific proteases. Depending on the source and preparation procedure of the samples, proteases and intra-molecular rearrangements, such as disulfide bonding, can vary. With for example 25 liquor cerebrospinalis (CSF) as sample source, the sequences "R-R" or "R-K" are frequently preceding a peptide's N-terminal position in a precursor as they are recognition sites of the prohormone convertase PC2 in CSF. Next to known enzyme recognition sites, some amino acids are more frequently and others are less frequent. Positions preceding or following N- and C-terminal positions of peptides can be predicted on the basis of their mere percentage 30 occurrence in a particular sample treated in that particular way. This kind of information can easily be determined empirically and an example for peptides present in human liquor cerebrospinalis is shown in the tables in Figures 14a to 14d. The tables "CSF: amino acid Before First Cleavage", "CSF: amino acid After First Cleavage", "CSF: amino acid Before Last Cleavage", and "CSF: amino acid After Last Cleavage" summarize empirically found N- 35 or C-terminal amino acid frequencies as a result of proteolytic processes. Rule h increases the fitness value when those amino acids at the top of the tables shown in Figures 14a to 14d are present at the corresponding positions in the predicted sequence, while those amino

acids at the bottom of these tables decrease the fitness value of the prediction. The tables shown in Figures 14a to 14d are suitable to predict the likelihood of the presence of certain amino acid residues at the N- or C-terminus of peptides present in human liquor cerebrospinalis as long as the CSF samples are treated in the same way as the CSF samples of the examples of the present invention. Tables similar to the ones shown in Figures 14a to 14d can be generated empirically for any sample such as whole blood, serum, plasma, urine and the like, and the treatment of the samples can be of any kind, as long as all samples are treated in the same way.

10 Rule h:

If the mass difference between the peptide coordinates of a known and an unknown peptide can be explained by the loss of one or more distinct N- or C-terminal amino acids, the fitness value of this prediction is increased.

15 Rule i:

If a prediction has been generated by one of the rules b to h or a combination thereof, proposing that the unknown peptide is a reactant or a product of a post-translational modification of the known peptide, this proposal is tested by determining in terms of accessibility of the reaction site within the protein sequence by an enzyme performing the given post-translational modification. Thus, if a look-up in a database storing three-dimensional data of peptides or proteins reveals that the proposed site is on the surface of the protein and/or its conformation sterically allows action of that enzyme, the fitness value of that prediction is increased. In the same way, if a region of a sequence is proposed to be modified by a post-translational modification process, the accessibility of that sequence 20 region to enzymes is assessed by means of algorithms estimating the hydrophobicity of that particular region (Engelman et al, *Ann. Rev. Biophys. Chem.*, 15, 321, 1986; von Heijne, *Eur. J. Biochem.*, 116, 419, 1981). For example, a highly hydrophilic sequence region is more likely to be accessible by enzymes performing post-translational modifications than a hydrophobic sequence region, thus the fitness value of that prediction is increased.

25

The results computed by applying rules a to i and optionally additional rules can be stored in a list or a database in computer readable format and/or can be printed or displayed via an appropriate user interface such as a monitor. If more than one prediction for an unknown amino acid sequence fits the results obtained with the rules described above, then the predicted sequence can be ranked with the best fitting sequence for the unknown peptide on top as shown at step 148 in Figure 8b. If the known peptide P1 has more than one related, unknown peptide P2, than the approach described can be repeated for all unknown peptides

P2 as shown in Figure 8c. The approach described above can be extended to any known peptide signal P1 in a list of peptides as exemplified in Figure 8d.

#### Differential Network Module

5 According to the present invention the interaction of the Differential Network Module with the fundamental CAN Module allows to identify peptides which independently from each other distinguish between a sample A and a sample B. A status can be young, old, healthy, diseased, sweet taste, bitter taste, transfected, non-transfected, yellow, green, male, female, pregnant, non pregnant, smoker, non smoker or any other criterion defining a group or a  
10 subgroup of samples or organisms from which samples are derived. Optionally the Differential Network Module is linked with various databases, containing data such as the status of the samples, as well as with the other modules of the present invention and especially the basic CAN Module as shown in Figure 1. The Differential Network Module instructs the CAN Module to define subgroups of samples defined by distinct criteria, such  
15 as the status of the samples, and further to calculate separately the peptide-to-peptide relations for any status or any combination of more than one status. First, those peptide pairs that suffice a threshold of correlation in a group of samples representing the status A, second, those peptide pairs that suffice a threshold of correlation in a group of samples representing the status B, and third, relations can be defined on the basis of differences  
20 between the correlations of the compared status A and status B. If a user is interested in peptide-to-peptide relations, that are most different in samples from two different statuses A and B, then he will search for peptides where the correlation coefficients of the respective peptide-to-peptide relations is different, and where  $\Delta r = |r_{\text{Status A}} - r_{\text{Status B}}|$  is preferably greater than 85% of all peptide-to-peptide  $\Delta r$ .

25 Figure 15 shows a flow chart schematizing the above described procedural steps of an interaction of the Differential Network Module with the CAN Module according to the present invention allowing for the identification of peptides suitable to be used as marker panels using the peptide topology of a plurality of samples taken from at least two different  
30 experimental groups representing a status A and a status B. At step 170 a respective mass spectrum for each sample of said plurality of samples is provided, wherein signal intensity peaks correspond to potential peptides. Then the measures of correlation between the signal intensities of said potential peptides are computed at step 172 for each plurality of samples within each experimental group separately. Finally pairs of potential peptides are selected at  
35 step 174, which exhibit a difference in the degree of correlation between the different experimental groups above a certain threshold, thereby providing peptides which are suitable

to be used as marker panels for diagnostic purposes to distinguish between status A and status B.

The results of the Differential Network Module allow statements about the different relations of peptides within samples of status A compared to status B as follows: If the difference of correlation coefficients of peptide I with peptide II in status A minus the corresponding correlation coefficient in status B is greater than a given threshold, signal coordinates of the peptide pairs, their mutual distance within the observed status A and status B or the degree of difference or combinations of the latter information are stored in a database or list. The Differential Network Module optionally provides the same visualization methods as the other modules, that means peptide coordinates and their relations can be represented as bullets connected by lines, respectively, as shown in Figure 6, and identified peptides can be reviewed via convenient connections to databases or lists containing supportive data resources.

Another use of this aspect of the present invention is the comparison of the molecular masses of peptides present in at least three samples, representing one or at least two different states, status A with corresponding samples and status B with corresponding samples. For example samples from individuals with a certain disease versus samples from individuals without that certain disease, samples from pregnant versus samples from non-pregnant individuals, samples from bacteria transformed with an expression vector versus samples from non-transformed bacteria, samples from yoghurt with a strong acidic taste versus samples from yoghurt with a mild acidic taste, etc. might be compared by computing the correlation measures of peptides present in these samples. The comparison of measurement parameters of a peptide within two samples corresponding to two different states A and B may also indicate that the peptide is present only in samples of state A but not in samples of state B. Also in this case the measurement parameters of this peptide in status A and status B possibly can be related by a measure of correlation. If at least two different peptides, e.g. peptide I and peptide II, are identified, the measurement values of the parameters for peptide I and peptide II can be combined. Using measurement values of at least three samples being representative of status A and three samples being representative of status B, a mathematical function can be computed. This mathematical function describes the correlation-network of peptide I and peptide II. It is possible to include more than two different peptides in one correlation-network, e.g. to include more than two different peptides in one mathematical function describing a correlation-network. The resulting mathematical function describes which combinations of measures of correlation of at least two peptides (peptide I and peptide II) allow to distinguish status A from status B.

Furthermore, another use of this aspect of the present invention comprises the automated identification of sets of peptides that allow a prediction of a status of a sample by a regression model. The invention detects relations between at least two peptides, where the relations are representative for a given status A. In a next step, a linear or non-linear regression model is set up that uses input parameters of the found peptides, such as their respective MALDI signal intensities, and that fits these input parameters to an end point parameter, such as the diagnosis (yes/no = 1/0), or that fits to another parameter of a peptide of this derived set.

5 10 In order to check whether a sample of unknown status is a member of the status A, the input parameters of these peptides from that sample are applied to the derived model. If the output value obtained from that sample deviates in the range as other samples from status A from an expected value obtained by means of the determined function, than this unknown sample can be considered to be from status A. Otherwise, the sample most likely has  
15 another status.

#### **Marker Panel Network Module**

According to the present invention the interaction of the Marker Panel Network Module with the fundamental CAN Module allows to identify peptides which independently from each  
20 other distinguish between a sample representing status A and a sample representing status B. For example a disease is caused by different factors such as inflammation and an increased heart beat rate. Each of these disease factors might result in altered concentrations of distinct peptides in for example blood plasma of the patient. If a panel of for example two peptide markers is used for diagnosis of the disease it would be useful if  
25 one of the peptide markers indicates inflammation and the other peptide marker indicates increased heart beat rate. The combination of these two markers would increase the specificity and sensitivity of the marker panel to detect the disease caused by a combination of inflammation and increased heart beat rate. The Marker Panel Network Module selects those potential peptides which are related to the disease but are most likely associated to  
30 different disease factors (in this hypothetical case inflammation and increased heart beat rate), since these peptide coordinates have a low measure of correlation to each other but both have a high correlation to the disease. Thus the specificity and sensitivity of a diagnostic test can be improved by combining these complementary peptide coordinates to a marker panel.

35

For example a disease 1 (status A) which is associated with inflammation has to be distinguished from another disease 2 (status B) which is not associated with inflammation.

There are, for example, four peptides found, which distinguish disease 1 from disease 2. Peptide 1 and peptide 2 are fragments from the same protein, for example from TNF-alpha, peptide 3 is, for example, a fragment of IL-6 and peptide 4 is a fragment of an unknown protein. All of these four peptides differentiate between disease 1 and disease 2 by a 5 measure of correlation, but peptide 1 and 2 correlate to each other, which is not surprising, as they originate from the same molecule (TNF-alpha). Additionally peptide 1 and peptide 3 correlate to each other, which is also not surprising, as TNF-alpha and IL-6 have similar pro-inflammation functions. Consequently there are two groups of peptides, peptides 1, 2 and 3 10 belong to one group and peptide 4 represents the second group. To obtain a diagnostic test, with improved specificity and/or sensitivity combination of the detection of peptide 1 and 2 or 1 and 3 or 2 and 3 would not increase the specificity and/or sensitivity as much as combination of peptide 1 and 4 or 2 and 4 or 3 and 4 would do. This method allows to identify panels of peptides with additive or synergistic value (diagnostic, therapeutic, functional, etc.).

15

Figure 16 shows a flow chart schematizing the procedural steps of an interaction of the Marker Panel Network Module with the CAN Module according to the present invention allowing for the identification of peptides suitable to be used as marker panels using the peptide topology of a plurality of samples taken from at least two different experimental 20 groups representing a status A and a status B. At step 180 a respective mass spectrum for each sample of said plurality of samples is provided, wherein signal intensity peaks correspond to potential peptides. Then potential peptides correlating with a parameter being representative of status A or status B are selected at step 182. Thereafter the measures of correlation between the signal intensities of said selected potential peptides for each plurality 25 of samples are computed at step 184 and finally pairs of potential peptides which exhibit no correlation of their respective signal intensities above a certain threshold are selected at step 186, thereby providing potential peptides which are suitable to be used as complementing peptides in a marker panel for diagnostic purposes to distinguish between status A and status B.

30

In other words, the Marker Panel Network Module selects potential peptides which correlate with a parameter being representative for status A or status B. The Marker Panel Network Module then queries the Correlation Associated Network (CAN) Module for those pairs of selected peptide coordinates, which have a very low measure of correlation of their 35 respective signal intensities to each other. The result are pairs of peptides which are related to the status A or B but not directly related to each other and can be combined for a marker

panel to distinguish between status A and B. It is possible to combine two or more peptides to a marker panel.

5 The Differential Network Module described in the previous section discovers combinations of peptides, whose ratio of concentration indicate a certain state and deviations from that ratio indicate a different state. It is mandatory to measure the signal intensity (e.g. concentration) of both/any peptide to calculate said ratio. The relations between two peptides may be present only in state A, whereas the relations between the same two peptides may be different or absent in state B.

10

In contrast, any peptide found by the Marker Panel Network Module described in the present section could serve as a diagnostic marker alone, but a combination of both markers improves the sensitivity/specificity etc. of the diagnostic test. The members of a marker panel ideally should not correlate with each other in any of both states. If the members of a 15 marker panel correlate with each other their combination most likely would not improve the sensitivity/specificity of the diagnosis.

### **Surrogate Network Module**

20 The Surrogate Network Module relates to the identification of peptides (so called surrogate peptides) that can replace or complement established diagnostic or therapeutic peptides or peptides of other use. If for instance it is discovered that peptides correlate with known bioactive therapeutic peptides, these peptides might serve as surrogates for therapeutic measures or even may exhibit a higher/larger potency, efficacy, specificity, selectivity and/or less undesirable side effects. These kind of peptides can be found using the Surrogate 25 Network Module in combination with the CAN Module according to the present invention by applying the steps shown in Figure 17. Initially a respective mass spectrum for each sample analyzed is provided, wherein signal intensity peaks correspond to potential peptides (step 190). Thereafter at step 192 measures of correlation between the signal intensity of a known peptide and the signal intensities of potential peptides are computed and finally those 30 potential peptides are selected at step 194, which exhibit a degree of correlation with the known peptide above a certain threshold, thereby providing potential peptides suitable to replace or complement the known peptide. Two exemplary applications of the Surrogate Network Module are given below

35 For example a plasma sample is known to contain the peptide insulin and a potentially unknown peptide X within the same plasma sample correlates with the peptide insulin. In this case peptide X might have the same function as insulin, as its correlation measure indicates

that it is related to insulin. The reason for this could be that peptide X is a derivative of insulin, for example a glycosylated form of insulin, or another peptide which is completely different from the amino acid sequence of insulin but which is involved in the same functional or metabolic cycles as insulin. In both cases peptide X could serve as an alternative to the 5 use of insulin for example in treating diabetes. It might also turn out that peptide X in combination with insulin improves the therapeutic effect of insulin by itself.

In a further example a tissue sample of a prostate cancer patient contains the prostate-specific antigen (PSA) peptide, which is a known marker for prostate cancer. Another 10 potentially unknown peptide Y is related by a correlation measure to the PSA peptide and consequently peptide Y might have the same diagnostic value as a biomarker for prostate cancer as the PSA peptide or the measurement of peptide Y might complement the prostate cancer diagnosis by PSA measurements.

15 **Interaction of Modules**

Though any of the modules described above can be used independently, any combination of these modules can be used and potentially can synergistically improve the result of one or more of the modules.

20 For example results of the Surrogate Network Module can be analyzed by the Sequence Network Module. In case the Surrogate Network Module yields peptide signals, which are not yet sequenced, a prediction of the sequence can give early hints for biological interpretation, thus accelerating validation processes of for example therapeutic or diagnostic peptides. However, a subsequent identification of these peptides by sequencing is recommended.

25 Results of the Differential Network Module can be analyzed with the Surrogate Network Module. If the Differential Network Module yields for example potential biomarkers, it is highly desirable to identify possible surrogate markers that show a similar behavior and therefore are of interest as well. Therefore a combination of the Surrogate Network Module 30 with the Differential Network Module accelerates the discovery of novel therapeutic, diagnostic or other peptides and is highly synergistic.

Furthermore, results of the Differential Network Module can be analyzed with the Sequence Network Module. If the Differential Network Module yields peptide signals, which have not 35 been sequenced yet, the prediction of the sequences of the unknown peptides can give early hints for biological interpretation, thus accelerating validation processes of potential

therapeutic, diagnostic or other peptides. However, the later identification of these peptides by sequencing is recommended.

### Examples

5 The following examples are intended to describe how the methods according to the present invention can be applied to real data. For the sake of a clarity only a limited number of exemplary measurement parameters are calculated and presented in the figures. However, as is readily observable by the person skilled in the art, the advantages of the methods according to the present invention become even more obvious when applied to large sets of  
10 data. On present computer systems commonly measures of correlation for data sets consisting of up to 6.000 potential peptides are commonly calculated and without undue effort data sets of up to 100.000 potential peptides can be analyzed by means of the methods according to the present invention.

15 **Example 1**

The basic CAN Module calculates to what extent a potential peptide for each individual potential peptide measured in a sample correlates to every other potential peptide in that sample. The CAN Module determines a network of correlations among the peptides which in case of some degree of correlation supposedly are related to each other for certain reasons  
20 such as a common precursor as the origin of the peptides or the same biological function of the different precursors of the correlating peptides.

In the present example the set of data, i.e. the data matrix, consists of 444.000 values comprising measurement parameters, in this case signal intensities, of 74 independent  
25 samples, each sample resulting in 6.000 peptide coordinates. The tables shown in Figures 18a, 18b list the corresponding raw data for four out of a total of 6.000 peptide coordinates. Four different methods to determine measures of correlation, namely, Spearman's rank order correlation, Pearson's product moment correlation, Kendall's rank correlation tau, and Minimal Spanning Tree (MST) diameter, are calculated for the three exemplified pairs of  
30 peptide coordinates comparing the peptide coordinate Fraction 54; m/z 2743.0 with three other peptide coordinates (Fraction 54; m/z 1371.5, Fraction 56; m/z 2927.2 and Fraction 20; m/z 1114.3) (see table shown in Figure 19). The definition of the threshold is an important step in the creation of Correlation Associated Networks and should be performed carefully as has been described in detail further above. In the data matrix  
35  $6.000 \times 6.000 \times 0.5 = 1.8 \times 10^7$  possible peptide-to-peptide pairs can be combined, and each of these pairs exhibits a certain correlation coefficient  $r$ . Figure 20 shows a plot of the probability of a peptide pair  $P(r)$  to have a certain correlation coefficient  $r$ . A value of  $r$  of zero

or close to zero describes relations which are completely random, whereas values of  $r$  close to 1 or -1 describes relations which respectively correlate or anti-correlate very strongly. The more peptide pairs are tested for a correlation by means of measures of correlations, such as Spearman's rank order correlation coefficient, the more peptide pairs by chance correlate 5 to some extent with each other. This means that a correlation coefficient regarded as informative and real has to pass a higher threshold value. It is recommended to perform a plot like in Figure 20 to estimate the information content of a given correlation coefficient. One curve (black circles) in this figure plots the likelihood (y-axis) for a given correlation coefficient (x-axis) for all peptide-to-peptide pairs from the said data matrix comprising 6.000 10 peptide coordinates. The other curve in Figure 20 marked by the white squares describes the likelihood of correlations occurring by chance.

Most probable true positive relations can be found where the area under the curve is small, while the maxima of the curves represent the correlations coefficients which are most likely 15 false positive relations. In case Spearman's rank order correlation coefficient is chosen as measure of correlation and  $|r_{threshold}| \geq 0.8$  is chosen as threshold for a definition of a peptide-to-peptide relation, the peptide coordinate Fraction 20; m/z 1114.3 is not related to the peptide coordinate Fraction 54; m/z 2743.0 (see table shown in Figure 19). In contrast, the peptide coordinates Fraction 54; m/z 1371.5, and F 56; m/z 2937.3 are highly related to the 20 peptide with coordinates Fraction 54; m/z 2743.0 (see table shown in Figure 19). These peptide relations could pass through a filter and be stored in a local Valentina Database file.

### Example 2

Assuming that one is interested in finding surrogate markers for Chromogranin A in 25 hypothetical prostate cancer patients and that some of the 74 samples described above originated from healthy male persons and some samples originated from prostate cancer patients. Under the further assumption that a peptide originating from Chromogranin A, amino acids 97-131, had been identified, the Surrogate Network Module would now query the basic CAN Module for peptide coordinates that are highly related by a correlation 30 measure with the hub-peptide Chromogranin A, 97-131. This could be done for example by defining that the Spearman's rank order correlation coefficient of peptide-to-peptide relations  $|r|$  has to comply with the relation  $|r| \geq 0.67$ . Then the Surrogate Network Module would instruct the CAN Module to query the Valentina Database, and report that there are about 14 peptide coordinates matching this condition. These peptide coordinates are searched in 35 databases for any known peptide fitting to these coordinates. In this way it would be found that three peptides known from the database and present in the list of 14 peptides belong to the Chromogranin/Secretogranin family as illustrated in the table shown in Figure 21. The

Surrogate Network Module would project the peptide coordinates of the related peptides and the hub peptide as bullets on a two-dimensional or three-dimensional plane, such as a peptide map fingerprint of a serum sample as shown in Figure 6. Relations between peptide coordinates are depicted as lines between the bullets. Lines can be selected by a computer pointing device such as a mouse and a small information window will pop up containing information about the kind of correlation measure and the value for the measure of correlation of the connected peptide coordinates is shown. The bullets can also be selected by a computer mouse click, and an information window will provide information about the peptide coordinate, and if this peptide coordinate has already been identified, then the name 5 of the precursor peptide, the start- and stop position of the identified peptide will be provided by retrieval of a "Sequence Information Database" as exemplified at 56 in Figure 1. Also links to other databases such as Swiss-Prot and GeneCard are provided and/or other databases such as the Patent database of the USPTO can be queried for the search terms 10 "name of the peptide" and "diagnostic". An internet browser window could display the results 15 from the US-Patent Database. The visualization of peptide-to-peptide relations and convenient connection and access to internet and intranet resources by the Surrogate Network Module significantly increases the speed of data acquisition that is needed for an evaluation of the results. The example of Chromogranin A indicates that other peptides originating from members of the secretogranin-chromogranin family are automatically found 20 by the CAN Module. These peptides are listed in the table shown in Figure 21 and can serve as a diagnostic marker for the prediction of the therapeutic success in the hypothetic prostate cancer patients.

### Example 3

25 In an exemplary hypothetic serum dataset 48 samples are derived from patients before prostatectomy and 26 samples from patients after prostatectomy. For the Differential Network Module a correlation measure, e.g. the Spearman's rank order correlation coefficient  $r$ , between the peptides is calculated for samples from patients before prostatectomy and for samples from patients after prostatectomy separately. The correlation 30 coefficient of Chromogranin A 97-131 and Secretogranin I 88-132 for all 74 samples is  $r = 0.67$ , for those patients before prostatectomy is  $r = 0.23$  and for those after prostatectomy is  $r = 0.97$  (see Figure 22). This shows that the peptides Chromogranin A 97-131 and Secretogranin I 88-132 obviously are much less related after prostatectomy than before. This also explains the loss of correlation for all patients. For the given example this means 35 that Secretogranin I 88-132 is a potential surrogate marker for Chromogranin A 97-131 only before prostatectomy, thereafter the relation is broken. This would have a significant impact on the design of a clinical evaluation of Secretogranin I 88-132 as a surrogate marker for

Chromogranin A, and could save enormous costs. Furthermore, the ratio of concentration of Chromogranin A 97-131 and Secretogranin I 88-132 is a diagnostic parameter itself. If the ratio deviates from 10/1 significantly, then a prostatectomy has been accomplished. Figure 22 exemplifies the use of the ratio of the signal intensities of Chromogranin A 97-131 and 5 Secretogranin I 88-132 as a diagnostic parameter: The ratio of 10/1 is present in all samples from patients before prostatectomy. In samples after prostatectomy this ratio is not present, i.e. the Secretogranin I/Chromogranin A relation is "broken".

#### Example 4

10 This section exemplifies the identification of representative peptides, also called "landmark peptides" and also refers to the given data matrix of 74 observations of 6.000 peptide coordinates already discussed in a previous example.

15 Two peptide coordinates are considered as related if the Spearman's rank order correlation of their signal intensities is above  $|r| > 0.8$ . The number of relations  $k$  a respective peptide has with different peptide coordinate is shown in the second row of the table shown in Figure 18a. From all peptide coordinates, Fraction 54; m/z 2743.0 has the most relations, i.e.  $k = 20$ . Therefore, this peptide coordinate would be No. 1 in a prioritization list. Then, the signal variance of Fraction 54; m/z 2743.0 is removed from the signal intensities of the 20 related 20 peptide coordinates, wherein Formulas 1, 2 and 3 are applied. Then the data of Fraction 54; m/z 2743.0 is removed from the data matrix. The tables shown in Figures 23a and 23b show the values given in the tables shown in Figures 18a and 18b after the variance of Fraction 54; m/z 2743.0 on the related peptide coordinates has been removed. This process is iterated to determine the next peptide coordinate as a candidate for the sequencing 25 prioritization list, until the number of peptides to be sequenced has been reached.

#### Example 5

30 In this example, the signal intensities of four fictive peptide coordinates of 74 samples, their respective mass-to-charge ratio and their fraction numbers are given (see table shown in Figure 18a). The calculation is performed using five fictive peptide coordinates using as the 5<sup>th</sup> peptide coordinate F 53; m/z 2823.0. One of the five signal coordinates, the fictive peptide HP 25-48 in Fraction 54; m/z 2743.029, has already been identified, and guided by the rules for the Sequence Network Module, the identities of the four remaining, unknown peptides will be proposed.

The measure of correlation of the four unknown peptide coordinates with HP 25-48 has been calculated in the CAN Module by means of Spearman's rank order correlation coefficient:

r (HP 25-48 and F 20; m/z 1114.3) = +0.00  
5 r (HP 25-48 and F 54; m/z 1371.5) = +0.92  
r (HP 25-48 and F 56; m/z 2927.3) = +0.84  
r (HP 25-48 and F 53; m/z 2823.0) = +0.87

As can be seen in Figures 24a to 24c and by the low correlation coefficient and MST 10 diameter shown below, respectively, F 20; m/z 1114.3 is not related to HP 25-48, thus will not be hypothesized to be related to the HP precursor protein. The generation of proposals for this peptide coordinate stop at this point.

In the same manner, the MST diameter was calculated as a measure of correlation:

15 MST diameter (HP 25-48 and F 20; m/z 1114.3) = 29 (see Figure 25a)  
MST diameter (HP 25-48 and F 54; m/z 1371.5) = 50 (see Figure 25b)  
MST diameter (HP 25-48 and F 56; m/z 2927.3) = 38  
MST diameter (HP 25-48 and F 53; m/z 2823.0) = 40 (see Figure 25c)

20 In contrast, peptide coordinates F 54; m/z 1371.5, F 53; m/z 2823.0 and F 56; m/z 2927.3 are highly related to HP 25-48 (see Figures 24b, 24c and Figures 25b, 25c). A proposal using the sequence of the precursor of the protein HP will be assigned to these peptide coordinates and the rules according to the Sequence Network Module of the present invention will be applied for sequence prediction.

25 Rule a determines whether the related peptide coordinate is a n-charged ion of HP 25-48. The calculation of Mass<sub>Deviation</sub> is exemplified with n = 1, 2, 3 or 4 and the mass-to charge ratios of F 54; m/z 1371.5 and F 56; m/z 2927.26 given in the table shown in Figure 26, using Formula 4. It is highly probable that F 56; m/z 1371.5 is a double charged ion of 30 HP 25-48, as in the case of n = 2 Mass<sub>Deviation</sub> < Mass<sub>Threshold</sub> = 0.5, therefore it is proposed as HP 25-48<sup>2+</sup>, i.e. the double charged ion of HP 25-48.

Rules b to i will now be applied to F 53; m/z 2823.0 and F 56; m/z 2927.3. Rule b assumes that the relation of the hub peptide P1 in fraction F 54; m/z 2743.029 of known identity with 35 the unknown peptide P2 (peptide coordinate F 53; m/z 2823.0) is derived from a post-translational modification. In this case, the mass difference of the hub peptide P1 and the

unknown peptide P2  $M_{\text{DIFF}} = |M_{\text{P1}} - M_{\text{P2}}| = 79.971$  might be caused by phosphorylation or sulphation (see table shown in Figure 9). Alignment of HP 25-48 with recognition sequence motifs of protein kinases, that are enzymes responsible for phosphorylation of proteins and peptides, identifies the sequence HP 35-37 to be "TYD", which is as a potential target of a 5 hypothetical protein kinase HPKC. Therefore a proposal for F 53; m/z 2823.0 is HP 25-48 with one phosphorylation at the tyrosine residue on position 36 of the peptide HP 25-48.

As stated before, if the unknown peptide and the known hub peptide are related, it is hypothesized that the unknown peptide is derived from the same precursor protein and thus 10 has the same precursor sequence as the known hub peptide. An algorithm systematically defines putative start and end positions, I an E, in the precursor sequence of the hub peptide P1 proposing a putative sequence fragment potentially derived from the precursor sequence, that could be the sequence of the unknown peptide P2 (see Figure 8b). Of course, the sum 15 of masses of the amino acid residues, plus their amino- and carboxy-terminal ends and plus potential posttranslational modifications must match the measured m/z ratio  $M_{\text{found}}$  of the unknown peptide P2 within a given threshold T. The Mass of the putative sequence is calculated by summing up the masses of the amino acid residues comprising the putative sequence for P2 plus the mass of a hydrogen and a hydroxyl group. Exemplary values of masses applying Formula 5 are given in the tables shown in Figures 11 and 12.

20 With HP 25-48 as the hub peptide and P2 having the peptide coordinate Fraction 56; m/z 2927.3 the Sequence Network Module searches for possible sets of start and end positions in the protein precursor sequence of HP as defined in Figure 27, that have a deviation of mass lower than the threshold  $T_{\text{MASS}} = 0.5$ .

25 One possible combination is a start position at amino acid No. 25 and an end position at amino acid No. 50 of HP resulting in the potential peptide HP 25-50:  
 $n_D = 2, n_A = 4, n_H = 2, n_K = 2, n_S = 1, n_E = 3, n_V = 2, n_R = 1, n_F = 1, n_L = 3, n_G = 1, n_I = 1, n_N = 1, n_T = 1, n_Y = 1$  in Formula 5 results in  $M_{\text{CALC}} = 2927.337$

30 This proposal is added to the list of proposals for P2.

The Sequence Network Module will now address the evaluation of the proposal HP 25-50 for P2 by applying rules c to i. In rule d, the chromatographic fraction of the proposed sequence  $F_{\text{CALC}}$  is estimated and compared with the found peptide coordinate of P2 ( $F_{\text{FOUND}}$ ). If  $F_{\text{CALC}}$  35 deviates from  $F_{\text{FOUND}}$  by less than the threshold for fractionation ( $T_{\text{FRACTION}}$ ) then the proposal is awarded with 2 model fitness points. If Formula 6, "Estimation of fraction number based on proposed sequence", is applied to HP 25-50, the calculated Fraction results in  $F_{\text{CALC}} = 56$ .

As P2 HP 25-50 is found in fraction 56, the number of model fitness points for this proposal is increased by two points. Formula 6 was generated empirically from a mathematical model using data originating from liquor cerebrospinalis samples separated using a specific HPLC-column (as described in the patent application WO 03/048775 A2) using a specific software.

5 Of course, for different types of samples and different separation methods other empirically determined models can be calculated in the same way.

Rule e rewards those proposals for P2, whose start-positions match the start positions of the hub peptide P1. In the case of HP 25-48 as the P1 hub peptide and HP 25-50 as the 10 proposal for the related peptide P2, the proposal HP 25-50 will be rewarded with 3 model fitness points.

Rule f rewards those proposals for P2, whose end-positions equal the end-position of the hub peptide P1. This is not the case with HP 25-50 as a proposal, therefore this rule does 15 not increase the model fitness points of this proposal for P2.

Rule g will increase the model fitness points of the proposal HP 25-50 by three points as the start position 25 is preceded by the amino acid sequence "R-R" (written in 1-letter amino acid code). The sequence "R-R" is a recognition site of prohormone convertases, which 20 commonly cleave after the second "R". In addition, rule g will increase the model fitness points for this proposal by another 3 points, as the "D-A" sequence is one of the preferred starts for peptide sequences present in liquor cerebrospinalis. Further sites of frequent proteolytic cleavage sites at start positions awarded by rule f are well known in the art.

25 Rule g assumes that the unknown peptide P2 is a product of N- or C-terminal proteolysis of the known hub peptide P1 or vice versa. The mass difference of P1 and P2  $M_{DIFF} = |M_{P1} - M_{P2}|$  is determined and aligned with the masses of the amino acids preceding and following the start- and end positions of P2 in the precursor sequence HP. In the example of HP 25-48 as P1 and HP 28-50 as P2 the mass difference is  $M_{DIFF} = 184.2$  and can be explained by the 30 amino acids "I-A" ( $M_I + M_A = 184.2$ ) which are following the end position of P1. Therefore P2 fits the model and the model fitness points for this proposal for P2 are increased by 3 points.

Obviously, rules c to i can be examined in any order, and rules can be left out for biological considerations, but still any combinations and any omissions of these rules are within the 35 scope of this invention.

The process described above can be repeated for all unknown peptides coordinates P2, which are related with HP 25-48.

#### Example 6

5 This example demonstrates the advantages associated with the methods according to the present invention by combining Correlation-Associated Peptide Networks with recognition of probable cleavage sites for peptidases and proteases in cerebrospinal fluid, resulting in a model able to predict the sequence of unknown peptides with high accuracy. On the basis of this approach, for instance the identification of peptide coordinates can be prioritized, and a  
10 rapid overview of the peptide content of a novel sample source can be obtained.

Cerebrospinal fluid (CSF) is in close contact with many parts of the brain. CSF aids to maintain a stable chemical environment for the central nervous system and is a route to remove products of brain metabolism. CSF distributes a multitude of biologically active  
15 substances within the central nervous system. It is acceptable to assume that CSF mirrors the physiological and pathophysiological status of the brain and, therefore, peptides from CSF represent a source of potential diagnostic and therapeutic target molecules.

Here the correlational behavior of peptides from CSF is analyzed, derived from the same  
20 protein precursor in more detail and correlational dependencies for the prediction of putative sequences of unknown peptides are exploited. If one assumes that a known peptide and an unknown peptide signal of a peptide-to-peptide pair might have a common protein precursor, the known protein precursor sequence is analyzed for proteolytic cleavage sites which could explain the generation of a signal with a mass corresponding to that of the unknown peptide.  
25 It will be shown that the combination of statistical analysis (CAN) and recognition of possible cleavage sites for peptidases and proteases in CSF results in a model with high predictive power for correct assignment of an unknown peptide signal to a protein precursor or even a sequence, thus reducing the number of peptides to be sequenced.

30 After approval by the local ethics committees, written informed consent was obtained from patients involved in this study. Human CSF was collected by lumbar puncture from neurological patients without cognitive impairment (n=39) and from patients suffering from dementia such as vascular dementia, Lewy-body dementia, frontotemporal dementia or Parkinson's disease (n=27). All CSF samples were prepared using mild conditions  
35 minimizing risk of sample alteration: The fluid was collected without aspiration and avoiding blood contamination. Samples were centrifuged for 10 min at 2000 g and the supernatant was stored at -80 °C until analysis.

Peptides were separated using reversed-phase C18 chromatography. 300 to 1500  $\mu$ L CSF was diluted 1:3.75 with water and the pH was adjusted to 2-3. Samples were loaded onto RP silica columns (250 x 4 mm column, Vydac, Hesperia, CA, USA; HP-ChemStation 1100 Agilent Technologies, Palo Alto, CA, USA). Retained peptides were eluted using an 5 acetonitrile gradient (4 to 80 %) in 0.05 % trifluoroacetic acid, collected into 96 fractions and lyophilized. Elution was monitored by UV detection. The retention time of major peptide peaks from repeatedly loaded extracts was used to confirm the reproducibility of the method.

After lyophilization, each HPLC fraction was resuspended in matrix solution (mixture of  $\alpha$ - 10 cyano-4-hydroxycinnamic acid and L-fucose (co-matrix) in 0.1 % acetonitrile/ trifluoroacetic acid (1:1 v/v) and applied to a matrix-assisted laser-desorption / ionization (MALDI) target, followed by ambient temperature air drying. Sample ionization was performed by application of repeated single laser shots over a representative area of the sample spot. The 15 accelerated ions were analyzed in a time-of-flight (ToF) mass spectrometer (Voyager-DE STR, Applied Biosystems, Framingham, MA, USA) in linear mode.

Peptides of interest were identified by mass-spectrometric sequencing using nanoESI-qTOF- 20 MS/MS (QSTAR pulsar, Sciex, Toronto, Canada) with subsequent protein database searching. The resulting peptide fragment spectra were acquired in product ion scan mode (spray voltage 950 V, collision energy 20-40 eV). Up to 200 scans per sample were 25 accumulated. Data processing previous to database searching included charge-state de- convolution (Bayesian reconstruct tool of the BioAnalyst program package, Sciex, Concord, Canada) and de-isotoping (customized Analyst QS macro; Sciex, Concord, Canada). The resulting spectra were saved in MASCOT (Matrix Science, London, UK) generic file format and submitted to the MASCOT search engine. Cascading searches including several 30 posttranslational modifications in Swiss-Prot (Version 39 or higher, [www.expasy.ch](http://www.expasy.ch)) and MSDB (Version 030212 or higher, EBI, Cambridge, UK) were performed by MASCOT DAEMON client (Version 1.9, Matrix Science) that allows, beside sequence determination, identification of modified amino acids as well as determination of their position within the peptide's sequence.

All mass spectra with the same fraction number of chromatography were baseline-corrected 35 averaged, and all 96 averaged mass spectra fractions were visualized in a "2D gel-like" format (peptide display), yielding an averaged peptide display (see Figure 28). Each peak (mass spectrometric signal) is depicted as a bar with its grey-scale intensity corresponding to the signal intensity of the corresponding MALDI-peak which corresponds to the relative amount of the peptide measured. The x-, y- and z-axis represent mass to charge ratios

(m/z), chromatographic fraction and mass spectrometric signal intensity, respectively. Masses range from 1,000 to 15,000 m/z ratios (x-axis). The data matrix of an individual peptide display consists of 16 million data points, of which, for a given signal to noise ratio, signal coordinates are extracted. For all samples, an identical set of signal coordinates is 5 thus present that is utilized for statistical analysis.

Data pre-processing of the acquired MALDI-ToF-mass spectra was performed applying baseline correction (RAZOR Library 4.0, Spectrum Square Associates, Ithaca, NY, USA) in combination with normalization of the mass spectra to a constant integral value. For the 10 benefit of simplicity and uniformity, all m/z-ratios were stated as average masses of the uncharged analyte. Wherever necessary, data was made available for the model by transformation of m/z-ratio data into this format.

For the analysis of all peptide-to-peptide relations, calculations of correlation were performed 15 with the signal intensities (i.e. relative peptide quantity) of all present (unknown) peptide coordinate data sets to any known peptide coordinate in peptide displays of patient samples: Any pair-wise relation of two peptides was rated by Spearman's rank order correlation of their respective signal intensities in all samples. Peptide pairs in combination with m/z ratios, chromatographic fraction and Spearman's rank order coefficients of correlation were stored in 20 a local Peptide-to-Peptide database.

In an automated approach, all peptide coordinates were individually queried in a peptide sequence database. The following rules were applied

- For each peptide coordinate:
  - Search for entry in peptide sequence database to match mass-to-charge ratio and chromatographic fraction within given thresholds.
  - If peptide coordinate found in peptide sequence database:
    - Retrieve information of peptide coordinate (sequence, average mass, name, precursor protein, start- stop position on precursor protein, precursor protein sequence).
    - Store information in an individual list of identifications.

All members of this list of identifications are now utilized as "hub peptides" for subsequent 30 correlational analysis.

- For each peptide coordinate (= hub peptide) with entry in peptide sequence database:
  - Create individual Correlation Associated Network by retrieving Peptide Coordinates with correlation coefficients above a given threshold in Peptide-to-Peptide database, thus becoming member of the CAN of the hub peptide.

- o For each member of a CAN:
  - If peptide coordinate of CAN member is not found in Peptide Sequence database:
    - Analyze protein precursor sequence of hub peptide for putative sequences that approximately match the mass-to-charge ratio of the CAN member peptide coordinate: create a list of putative sequences that match the mass-to-charge ratio of the CAN member peptide coordinate within the range of mass accuracy (here: less than 500 ppm) by permuting the start- and end positions on the protein precursor sequence and simultaneously summation of masses of the amino acid residues of the putative sequences.
    - For each putative sequence in generated list:
      - o Evaluate putative Sequence (rules see below) with bonus points.
      - Rank putative sequences according to the number of bonus points.
  - o For each peptide coordinate
    - Present Top 3 putative sequences.

20 The determination of the bonus points is explained further below.

a) If the amino acid residue before/after amino-terminal/carboxy-terminal cleavage site of the putative peptide sequence on the precursor sequence corresponds to the following amino acid residues (one letter code), the proposal is awarded with the respective bonus points (bpt):

25	1. before amino-terminal cleavage site:	M : 2 bpt;	R: 5 bpt
	2. after amino-terminal cleavage site (N+1):	D : 3 bpt;	M: 2 bpt
	3. before carboxy-terminal cleavage site (C-1):	no rules	
	4. after carboxy-terminal cleavage site (C+1):	K : 3 bpt;	R: 4 bpt

30 b) If the pairs of amino acids before/after amino-terminal/carboxy-terminal cleavage site of the putative peptide sequence on the precursor sequence correspond to the following amino acid pairs, the proposal is awarded with the respective bonus points (bpt):

1. before amino-terminal cleavage site:	KR: 18 bpt;	RR: 22 bpt;
2. after amino-terminal cleavage site:	DA: 43 bpt;	GR: 11 bpt;
3. before carboxy-terminal cleavage site:	GA: 20 bpt;	QK: 20 bpt;
	VN: 16 bpt	
4. after carboxy-terminal cleavage site:	KR: 22 bpt	

c) If the putative sequence has the same start position as the known hub peptide, the proposal of this sequence is awarded with 69 bonus points. If the putative sequence has the same end position as the known hub peptide, the proposal of this sequence was awarded with 63 bonus points.

5 The determination of the bonus points is explained further below.

The peptidome of 66 independent CSF samples was analyzed using a combination of chromatographic separation (96 fractions) and subsequent mass spectrometry, leading to a database containing 7104 MALDI-ToF pre-processed mass spectra. All mass spectra with 10 the same fraction number were averaged, yielding an averaged peptide display (see Figure 28). 139 different peptides from 31 protein precursors have previously been identified by sequence determination from CSF-preparations carried out in an identical manner to the 66 CSF samples. The peptide coordinates were located on the averaged peptide display. Since abundant peptides can be found in more than one fraction, 224 instead of 139 peptide 15 coordinates were located. The MALDI mass spectrometric signal intensities of the 224 peptides coordinates were determined in each of the 66 samples. Spearman's rank correlational analysis was performed for any given signal-to-signal combination resulting in  $224^2/2$  correlations of peak signal intensities.

20 As described in detail further above, a network is defined as a collective of a peptide of interest, the so-called hub peptide, and peptides highly correlating with this peptide, selected from all peptides by exceeding an arbitrarily defined correlational threshold. This concept is exemplified by two networks of VGF and albumin peptides: The network of VGF 26-58 as a hub peptide (see Figure 29), calculated with a threshold correlation of  $|r| \geq 0.68$ , groups 25 peptides derived from several regions of the VGF precursor protein (see Figure 30). The network with the hub peptide albumin 25-48 (see Figure 31), calculated with a threshold of correlation of  $|r| \geq 0.67$ , predominantly contains peptides derived from the amino-terminal region of the albumin protein precursor. The observation of networks of VGF 26-58 and albumin 25-48 gave rise to the hypothesis, that a network predominantly contains peptides 30 derived from the same protein precursor, if the threshold of correlation is sufficiently high. This hypothesis was tested as follows: Networks of the 224 signal coordinates were created with increasing thresholds of correlation coefficients (see Figure 32). Members of these networks were predicted as derivatives from the same precursor protein, containing the hub peptide and the predictions matched with the peptides previously identified by ESI-MS/MS.

35 The table of Figure 32 lists the number of correct and false predictions of protein precursors: Predictions with a high threshold of correlation yield few predictions of precursor proteins; predictions with low threshold of correlation coefficients yield several hundred predictions of

precursor proteins. The correctness of precursor predictions reaches 100% at  $|r| \geq 0.95$  and rapidly decreases with lower thresholds of correlation.

It was assumed that any member of a network is derived from the same protein precursor  
5 without experimental determination of the sequence of the unknown signal coordinate. Thus,  
start- and end positions on that precursor protein sequence were systematically permuted  
and iterated resulting in putative peptide sequences. Figure 33 gives an overview of the  
nomenclature for start and end of a peptide on a protein precursor. Any combination of start-  
/end- positions resulting in a putative sequence that matches the m/z-ratio of the unknown  
10 peptide coordinate within a given mass tolerance of less than 500 ppm was considered as a  
valid proposal. Combinations that do not fulfill this mandatory criterion were rejected. The  
following investigations showed that some valid proposals were more probable than others,  
since their start/end positions are very likely sites of cleavage in human CSF. If a proposal  
meets one or more of the following criteria, its bonus points are increased.

15 Peptides are generated by cleavage of peptide bounds by proteolytic enzymes. These  
proteases recognize specific sites (amino acid sequence motifs), where a cleavage occurs.  
The probability of cleavage as a function of a particular amino acid and the position of the  
amino acid regarding the cleavage site was investigated and compared with the occurrence  
20 of the respective amino acid at any position in all precursor sequences. The table in Figure  
34 shows the results obtained from the data set enrolling the 139 peptides. For example, in  
31 % of all peptides, an arginine residue (R) preceeded the amino-terminal cleavage site,  
but the average content of arginine of the observed precursor sequences was only 6 %. It  
was concluded, that the probability of cleavage was 5 times increased if an arginine residue  
25 was found at the N-1 position. Rules were defined considering the different amino acids at  
the position of interest: an x-fold increase of probability of cleavage was awarded with x  
bonus points: For example, the bonus point score of an proposal with arginine at N-1 were  
increased by 5 points, reflecting the 5-fold increase of probability. The number of rules was  
30 limited in order to avoid excessive overfitting of the model: A rule had to be based on at least  
five peptides and a twofold increase of probability.

Besides investigations of single amino acids, pairs of amino acids were investigated for their  
influence on an increased probability of cleavage. The probability of cleavage as a function  
35 of a pair of amino acids and the position of such a pair with respect to the cleavage site was  
investigated and compared with the occurrence of the respective pair of amino acids at any  
position in all precursor sequences (see table in Figure 35). For example, arginine-arginine  
residues (RR) were found in 18 / 139 = 12.9 % of the peptides before the amino-terminal

cleavage site, while RR was found only at  $398 / 68516 = 0.58\%$  at any other position. Thus, the probability of cleavage after RR is  $12.9\% / 0.58\% = 22$ -fold increased compared to arbitrary positions. Consequently, a rule considering RR before amino-terminal cleavage increases the bonus point score of a corresponding proposal by 22 bonus points. Thus rules 5 were defined considering pairs of amino acid at the position of interests: an x-fold increase of probability was awarded with x bonus points. Yet, a rule had to be based on at least five peptides and a 10-fold increase of probability, otherwise it was rejected.

Many related peptides were found to have the same start position on a precursor protein as exemplified by VGF 26-58, VGF 26-59, VGF 26-61 and VGF 26-62 (see Figures 29 & 30 #1, 10 #4, #5, #8, #9) or Albumin 25-48, Albumin 25-45 and Albumin 25-50 (see Figure 31 #1, #4/#5, #2, #8/#9). Likewise, many sequences of related peptides stop with the same end 15 position: VGF 25-59 and VGF 26-59 (see Figures 29 & 30 #4/#5, #7) or Albumin 27-50 (see Figure 31 #6/#7) and Albumin 25-50 (see Figure 31 #8/#9). It was found that 14.1 % of all peptides from the same precursor in the data set had the same start position and 12.7 % the 20 same end position. The probability of two peptides having the same start- or end-position by chance was assumed as  $1/n$ , where n is the length of the precursor. With an average precursor length of  $n=492$  in the data set, the increase of probability was 69-fold for two peptides having the same start-position and 63-fold for two peptides having the same end position. Consequently, proposals with the same start-position as the hub peptide were 25 awarded with 69 bonus points and those with the same end-position with 63 bonus points.

The application of the above described rules is exemplified by two proposals (see table in Figure 36), verified by ESI-MS/MS: According to the model, the known peptide VGF 26-58 25 predicts the unknown peptide coordinate with  $m/z_{\text{average}}$  3688.0 as VGF 26-62, because the calculated  $m/z$ -ratio of the putative sequence matches the found  $m/z$ -ratio (prerequisite condition). Hub peptide and putative sequence have the same start position (+69 bonus 30 points), and the putative sequence terminates before an arginine residue (+4 bonus points) (see Figure 30 #1, #4). Moreover, VGF 26-58 predicts the unknown peptide with  $m/z_{\text{average}}$  2419.41 as VGF 350-370. The calculated  $m/z$ -ratio of the putative sequence matches the found  $m/z$ -ratio (prerequisite condition), a single arginine precedes (+ 5 bonus points) and a dibasic site RR follows the putative sequence (+22 bonus points) (see Figure 30 #1, #3).

In order to asses the power of prediction of the described model, 139 peptides identified by 35 ESI-MS/MS were split into a group of 70 peptides, that were used for the prediction of a second group of 69 peptides, whose sequence identity was suppressed during the prediction process (see Figure 37): For all 224 signal coordinates, the Correlation Associated Networks

were calculated with a coefficient of correlation of  $|r| = 0.75$ . Seventy signal coordinates were used to predict the precursor protein, start- and stop position and thus the sequence of signal coordinates, which were members of the network of the corresponding predicting hub peptide. For any signal coordinate, up to three proposals were listed, having the proposal with most bonus points on top. Six models with increasing levels of complexity were compared regarding their power in prediction of precursor protein, start and end positions (see table in Figure 38). Predictions for the second group of 69 peptides were compared with identifications by ESI-MS/MS. The statistics were distinguished for any stored proposal and for the proposal with most bonus points: In all models, the percentages of correctly predicted start and end positions is better for the proposal with most bonus points compared to those of all proposals. The percentage of correct precursor sequence and start-/end-position increased with complexity of the model. Model 5, combining two sets of rules considering both the probability of cleavage near single amino acids and pairs of amino acids, yielded better results than model 2 and model 3, which applied only one set of rules. Best results were achieved by model 6, which incorporated all described rules, in particular those proposals with most bonus points: 85 % of all proposals were correct regarding the precursor sequence, the start- and stop positions, 89 % of all proposals yielded a correct protein precursor prediction, and only 11 % were wrong in both.

The above example thus demonstrates that related peptides are automatically grouped by CANs. The underlying algorithm exploits the fact that concentrations of peptides from different steps of the processing chain can display a conserved ratio, as shown previously for CSF-derived peptides. These conserved ratios of related peptides were reliably found by Spearman's rank order correlation analysis, which is the basis for the definition of CAN relations. The results show that CANs can be used to automatically group intermediate products of peptide processing. At high thresholds of correlation coefficient, the number of predictions is low, but each having a high degree of accuracy. Decreasing thresholds delivers a growing number of predictions with false ones, finally outbalancing the correct predictions. The present example was based on the assumption/condition that a strict threshold delivers a network whose members are solely derived from the same protein precursor. This was the fundamental basis for the prediction of the sequence of unknown network members. Since the network was based on mass spectrometric data, all peptide signals were characterized by their mass-to-charge ratio. By iterating start and downstream end position on the protein precursor sequence of the previously sequenced hub peptide, putative peptide sequences were generated that matched experimental molecular masses of selected unknown MALDI mass spectrometric peptide signals. The mass accuracy of less than 500 ppm of the MALDI-ToF measurement in linear mode is sufficient to reduce the overwhelming number of theoretical combinations of start- and stop positions of a precursor to a concise selection. The putative sequences were evaluated by models based on the

presumed protein precursor's sequence in combination with found proteolytic cleavage patterns in human CSF. In this approach, posttranslational modifications were not considered, considerably reducing the degree of freedom for possible predictions. However posttranslational, as well as other modifications in general can be used to search for 5 correlation in peptide signals.

Six models were tested that were built on different sets of rules and combinations thereof. Since the cleavage of protein precursors is sequence- and tissue-specific, the sequence 10 specificity of proteases in human CSF was investigated: Pairs of amino acids, the 'motifs', at the four positions, before and after amino-terminal and carboxy-terminal cleavage site, were differentiated for cleavage pattern analysis (Figure 38).

The described rules applied to sequence prediction are generic since they are based on the x-fold increase of the probability for the given event, and scoring the respective proposals 15 with x bonus points. Combining the individual rules by summing the bonus points substantially increased the accuracy of prediction. This confirms that the rules of the different approaches are complementary and not contradictive. The magnitude of the bonus points was significantly different due the individual definition for every single parameter.

20 However, if the algorithm is applied to other sample matrices, the presented rules most likely will have to be redefined. The rules can be determined empirically using these other sample matrices. It is also recommended to test the parameters r and bonus points in a given data set with known peptides to determine the false positive rate prior to use for prediction of peptide sequences. The parameters should be readjusted until the false positive rate and 25 prediction number match the design and requirements of the experimental purposes.

As a result of the combination of statistical analysis and peptide biology for the definition of a set of specific rules, a promising model is envisaged that predicts the sequence of peptides with high accuracy. A system of bonus points was used to select the prediction, which fitted 30 the model best. Proposals with the highest score of bonus points were compared with ESI-MS/MS identifications, and were found to be correctly predicting 85 % of protein precursors, start- and end positions and 89 % of precursor protein only. A further improvement is expected by using MALDI-ToF measurements in the reflectron mode with a mass accuracy of less than 30 ppm, and with broader sequence coverage to redefine the model.

As a consequence of the promising results of this proof-of-concept study the following procedure is suggested, in case a rapid overview of a peptide content of a new sample source has to be obtained: The peptide coordinates of a novel samples source are defined based on a representative peptide display. Thereafter, the related peptide coordinates are  
5 determined by calculating the CAN of any peptide coordinate: Hubs with the most network members are considered as related to a multitude of other peptide coordinates, thus being the most representative ones (Lamerz et al., 2005, Proteomics, 5:x-xx). These peptide coordinates should be identified first. On the basis of these identifications, CAN is used to predict sequence of the remaining not identified peptide coordinates. The identification of  
10 peptide signal coordinates that are adequately described by the model can be postponed or discarded from the identification list, leaving more resources for the identification of less abundant peptides or peptides unsatisfactorily described by the model. This procedure can be repeated several times with the additional sequence information generated during the process, resulting in a reduction of MS/MS identification work while achieving a comparably  
15 deep insight in the content of the novel sample source.

Posttranslational and other modifications of peptide sequences can be included by scanning the sequence of the hub peptide for specific motifs characteristic or susceptible for such modifications, such as phosphorylation, dephosphorylation, oxidation, reduction,  
20 glycosylation, deglycosilation, acetylation and other modifications known for peptides. Subsequently, the mass difference between the hub peptide and its related peptides can be analyzed to assess whether the mass difference corresponds to the respective posttranslational modification. This implies the implementation of many, even thousands of motifs, as documented in PROSITE (Falquet et al., 2002, Nucleic Acids Res., 30:235-238),  
25 and the scanning process can be computationally laborious.

### Example 7

CANs were also exploited for the discovery of surrogates of biomarkers. The whole albumin molecule is routinely used in diagnosis as gold standard to determine the integrity of the  
30 brain barriers (blood-brain barrier, blood-CSF barrier). The ratio of albumin concentration in CSF and blood, the 'albumin ratio', correlates with the extent of a barrier disruption (Reiber et al., 1980, J. Neurobiol., 224:89-99), leading to the increased transfer of "blood born" peptides and proteins into the CSF. Previous work shows, that the albumin peptide representing amino acids 25-48 of human albumin can serve as a marker for an impaired  
35 brain barrier (Heine et al., 2002, J. Chromatogr. B Analyt. Technol. Biomed. Life. Sci. 782(1-2):353-61). By way of this example (see Figures 39 and 40) it is demonstrated that members of the CAN calculated by use of the albumin 25-48 peptide as hub peptide have the same potential as surrogate markers to assess disruption of the brain barrier as the albumin 25-48

peptide and thus that CANs are suitable to identify surrogate markers of known markers (in this case of albumin 25-48).

This was tested in an independent experimental set-up using well-documented CSF samples taken from patients with different severe disruptions of the blood-CSF barrier. Subsequently to the identification of these potential surrogates for albumin 25-48, the surrogates identified were searched in the original dataset described in a previous work (Heine et al., 2002, J. Chromatogr. B Analyt. Technol. Biomed. Life. Sci. 782(1-2):353-61). It was confirmed, that the surrogates proposed (see Figure 39 and 40) are suitable to diagnose the impaired brain barrier of the patients analyzed in a previous work (Heine et al., 2002, J. Chromatogr. B Analyt. Technol. Biomed. Life. Sci. 782(1-2):353-61). This shows, that CANs are suitable to identify surrogates of known markers, as claimed in this invention.

The samples used for the study performed with the present invention were e.g. human CSF collected by lumbar puncture from 74 patients suffering from vascular dementia, Lewy-body dementia, frontotemporal dementia, Parkinson's disease, depression, lumbago, facial paresis, vertigo, polyneuropathy or optic neuritis.

These samples were analysed by reversed-phase chromatography and MALDI mass spectrometry under the same conditions as the samples in example 6. Albumin 25-48 as the hub peptide displayed strong correlations ( $|r| > 0.75$ ) with 25 different peptide signals and, most importantly, a significant correlation with the albumin ratio in the new sample set, which was determined using standard albumin-ELISA tests, as known in the art ( $|r| = 0.73$ ). It was found that all network members correlated positively with the albumin quotient, and 16 out of 25 reached a significant level ( $|r| > 0.7$ ,  $n=9$ ,  $p<0.05$ ). This positive correlation with the established and accepted albumin ratio as a measure for blood-CSF barrier disruption indicate the correctness of the predicted peptide-to-peptide relations in CSF. Five prominent network members were identified subsequently as structurally similar amino-terminal fragments of albumin, namely albumin 25-48, albumin 25-50, albumin 25-51 and albumin 27-50 by sequencing. The novel peptide alpha-1-antitrypsin 397-418 of the albumin 25-48 CAN correlated even stronger to the albumin quotient ( $|r| = 0.83$ ) than the albumin fragment itself ( $|r| = 0.73$ ). The identification of the alpha 1-antitrypsin 397-418 peptide as a member of the albumin CAN highlights the power of the claimed methodology for the identification of new chemically unrelated peptide surrogates with a high diagnostic potential. Interestingly, alpha 1-antitrypsin as whole protein is already described as a protein directly correlating with disturbances in the blood-brain barrier determined by assessment of the ratio of albumin in

CSF to that in serum (Pearl et al., 1985, Arch. Neurolo. 42:775-777) further supporting that CANs are suitable to predict surrogates of known markers.

The person skilled in the art will appreciate from the foregoing that the range of application  
5 of CANs is expandable to any proteomic approach that allows a semi-quantitative analysis of components, for example data from two-dimensional gels (2D-gels). In such cases pair wise correlation coefficients of the components can be calculated, but it is of utmost importance to verify spot identity, avoiding inclusion of spots deriving from contaminating proteins. There,  
10 the precision of the two dimensions of peptidomics CAN, i.e. chromatographic fraction in RP-HPLC (usually better than 1 %) and in MALDI-MS (usually better than 100 ppm) is highly superior to separations obtained by 2D-gel electrophoresis (Schulz-Knappe et al., 2001, Comb. Chem. High Throughput. Screen., 4:207-217). On the other hand, CANs based on the approach described in the examples of the present invention are restricted to proteins < 15 kDa, while CAN based on 2D-gels can also address networks of larger proteins.

15 CANs are also applicable to peptide and protein quantification data from Isotope-Coded Affinity Tag (ICAT) mass spectrometric experiments. In ICAT experiments peptides and proteins present in the samples are isotopically labelled through a reactive group that specifically binds to cysteine residues. In a low molecular weight (peptidome) region, the  
20 number of peptides and small proteins that do not contain cysteins needed for ICAT labelling is higher compared to the proteomics field, thus decreasing the efficiency of ICAT. Novel labels such as iTRAQ which is an amine specific isotope labelling technique developed by Applied Biosystems, Foster City, CA, USA, will allow detection of all small proteins/peptides in CAN experiments.

25 It is envisaged that CANs will also support the interpretation of data from tryptic digestions of protein or peptide containing samples. Although the CAN methods presented here are based on undigested, native peptides, a similar clustering of different peptide species derived from the same precursor after tryptic digestion is possible. While this invention has been  
30 described with reference to preferred embodiments, it will be understood by those skilled in the art that various changes or modifications in form and detail may be made without departing from the scope of the invention as defined in the following claims.

For example it is readily apparent that the present invention can advantageously be utilized  
35 basically with all kinds of samples potentially containing peptides, such as samples from animals, plants, fungi, humans, parasites, microorganisms, such as bacteria, yeasts, viruses, and the like, samples from food or other agricultural materials such as meat, milk, grain,

vegetables, wool, cotton, silk, samples from cosmetic products or other products containing peptides such as cleaning agents (often containing proteolytic enzymes), etc. Samples for example can be plasma, serum, hemo-filtrate, whole blood, blood cells, tissues samples, in vitro grown cells, cell culture supernatants, urine, cerebrospinal fluid, lymph fluid, sputum, 5 tear fluid, ascites, preparations of cell organelles, tissue homogenate or homogenates of a virus, a microorganism, a parasite, a multi-cellular organism, an animal, a fungus or a plant and the like or combinations thereof. Examples of combinations are in vitro cultured cells infected with a microorganism or treated with pharmaceutical substances, tissue samples of humans infected with a microorganism, food products containing microorganisms, tissue 10 culture supernatants of cells treated with peptides or mixtures of peptides present in food or cosmetic products, and the like.